



Institute for Empirical Research in Economics  
University of Zurich

Working Paper Series  
ISSN 1424-0459

---

Working Paper No. 515

**Nonlinear Shrinkage Estimation of  
Large-Dimensional Covariance Matrices**

Olivier Ledoit and Michael Wolf

Revised version, December 2011

---

# Nonlinear Shrinkage Estimation of Large-Dimensional Covariance Matrices

Olivier Ledoit  
Department of Economics  
University of Zurich  
CH-8032 Zurich, Switzerland  
olivier.ledoit@econ.uzh.ch

Michael Wolf\*  
Department of Economics  
University of Zurich  
CH-8032 Zurich, Switzerland  
michael.wolf@econ.uzh.ch

October 2010

Revised December 2011

## Abstract

Many statistical applications require an estimate of a covariance matrix and/or its inverse. When the matrix dimension is large compared to the sample size, which happens frequently, the sample covariance matrix is known to perform poorly and may suffer from ill-conditioning. There already exists an extensive literature concerning improved estimators in such situations. In the absence of further knowledge about the structure of the true covariance matrix, the most successful approach so far, arguably, has been shrinkage estimation. Shrinking the sample covariance matrix to a multiple of the identity, by taking a weighted average of the two, turns out to be equivalent to linearly shrinking the sample eigenvalues to their grand mean, while retaining the sample eigenvectors. Our paper extends this approach by considering nonlinear transformations of the sample eigenvalues. We show how to construct an estimator that is asymptotically equivalent to an oracle estimator suggested in previous work. As demonstrated in extensive Monte Carlo simulations, the resulting *bona fide* estimator can result in sizeable improvements over the sample covariance matrix and also over linear shrinkage.

KEY WORDS: Large-dimensional asymptotics, nonlinear shrinkage, rotation equivariance.

JEL CLASSIFICATION NOS: C13.

---

\*Research has been supported by the NCCR Finrisk project “New Methods in Theoretical and Empirical Asset Pricing”.

# 1 Introduction

Many statistical applications require an estimate of a covariance matrix and/or of its inverse when the matrix dimension,  $p$ , is large compared to the sample size,  $n$ . It is well-known that in such situations, the usual estimator — the sample covariance matrix — performs poorly. It tends to be far from the population covariance matrix and ill-conditioned. The goal then becomes to find estimators that outperform the sample covariance matrix, both in finite samples and asymptotically. For the purposes of asymptotic analyses, to reflect the fact that  $p$  is large compared to  $n$ , one has to employ large-dimensional asymptotics where  $p$  is allowed to go to infinity together with  $n$ . In contrast, standard asymptotics would assume that  $p$  remains fixed while  $n$  tends to infinity.

One way to come up with improved estimators is to incorporate additional knowledge in the estimation process, such as sparseness, a graph model, or a factor model; for example, see Bickel and Levina (2008), Rohde and Tsybakov (2010), Cai and Zhou (2012), Ravikumar et al. (2008), Rajaratnam et al. (2008), Khare and Rajaratnam (2011), Fan et al. (2008), and the references therein.

However, not always is such additional knowledge available or trustworthy. In this general case, it is reasonable to require that covariance matrix estimators be rotation-equivariant. This means that rotating the data by some orthogonal matrix rotates the estimator in exactly the same way. In terms of the well-known decomposition of a matrix into eigenvectors and eigenvalues, an estimator is rotation-equivariant if and only if it has the same eigenvectors as the sample covariance matrix. Therefore, it can only differentiate itself by its eigenvalues.

Ledoit and Wolf (2004) demonstrate that the largest sample eigenvalues are systematically biased upwards, and the smallest ones downwards. It is advantageous to correct this bias by pulling down the largest eigenvalues and pushing up the smallest ones, towards the grand mean of all sample eigenvalues. This is an application of the general shrinkage principle, going back to Stein (1956). Working under large-dimensional asymptotics, Ledoit and Wolf (2004) derive the optimal *linear* shrinkage formula (when the loss is defined as the Frobenius norm of the difference between the estimator and the true covariance matrix). The same shrinkage intensity is applied to all sample eigenvalues, regardless of their positions. For example, if the linear shrinkage intensity is 0.5, then every sample eigenvalue is moved half-way towards the grand mean of all sample eigenvalues. Ledoit and Wolf (2004) both derive asymptotic optimality properties of the resulting estimator of the covariance matrix and demonstrate that it has desirable finite-sample properties via simulation studies.

A cursory glance at the Marčenko and Pastur (1967) equation, which governs the relationship between sample and population eigenvalues under large-dimensional asymptotics, shows that linear shrinkage is the first-order approximation to a fundamentally nonlinear problem. How good is this approximation? Ledoit and Wolf (2004) are very clear about this. Depending on the situation at hand, the improvement over the sample covariance matrix can either be gigantic or minuscule. When  $p/n$  is large and/or the population eigenvalues are close to one another, linear shrinkage captures most of the potential improvement over the sample covari-



depends on the population covariance matrix, which is unobservable, it is always safer *a priori* to use nonlinear shrinkage.

Many statistical applications require an estimate of the precision matrix, which is the inverse of the covariance matrix, instead of (or in addition to) an estimate of the covariance matrix itself. Of course, one possibility is to simply take the inverse of the nonlinear shrinkage estimate of the covariance matrix itself. However, this would be *ad hoc*. The superior approach is to estimate the inverse covariance matrix directly by nonlinearly shrinking the inverses of the sample eigenvalues. This gives quite different and markedly better results. We provide a detailed, in-depth solution for this important problem as well.

The remainder of the paper is organized as follows. Section 2 defines our framework for large-dimensional asymptotics and reviews some fundamental results from the corresponding literature. Section 3 presents the oracle shrinkage estimator which motivates our *bona fide* nonlinear shrinkage estimator. Sections 4 and 5 show that the *bona fide* estimator is consistent for the oracle estimator. Section 6 examines finite-sample behavior via Monte Carlo simulations. Finally, Section 7 concludes. All mathematical proofs as well as some further Monte Carlo simulations are collected in two appendices.

## 2 Large-Dimensional Asymptotics

### 2.1 Basic Framework

Let  $n$  denote the sample size and  $p \equiv p(n)$  the number of variables, with  $p/n \rightarrow c \in (0, 1)$  as  $n \rightarrow \infty$ . This framework is known as large-dimensional asymptotics. The restriction to the case  $c < 1$  that we make here somewhat simplifies certain mathematical results as well as the implementation of our routines in software. The case  $c > 1$ , where the sample covariance matrix is singular, could be handled by similar methods, but is left to future research.

The following set of assumptions will be maintained throughout the paper.

- (A1) The population covariance matrix  $\Sigma_n$  is a nonrandom  $p$ -dimensional positive definite matrix.
- (A2) Let  $X_n$  be an  $n \times p$  matrix of real independent and identically distributed (i.i.d.) random variables with zero mean and unit variance. One only observes  $Y_n \equiv X_n \Sigma_n^{1/2}$ , so neither  $X_n$  nor  $\Sigma_n$  are observed on their own.
- (A3) Let  $((\tau_{n,1}, \dots, \tau_{n,p}); (v_{n,1}, \dots, v_{n,p}))$  denote a system of eigenvalues and eigenvectors of  $\Sigma_n$ . The empirical distribution function (e.d.f.) of the population eigenvalues is defined as,  $\forall t \in \mathbb{R}$ ,  $H_n(t) \equiv p^{-1} \sum_{i=1}^p \mathbb{1}_{[\tau_{n,i}, +\infty)}(t)$ , where  $\mathbb{1}$  denotes the indicator function of a set. We assume  $H_n(t)$  converges to some limit  $H(t)$  at all points of continuity of  $H$ .
- (A4)  $\text{Supp}(H)$ , the support of  $H$ , is the union of a finite number of closed intervals, bounded away from zero and infinity. Furthermore, there exists a compact interval in  $(0, +\infty)$  that contains  $\text{Supp}(H_n)$  for all  $n$  large enough.

Let  $((\lambda_{n,1}, \dots, \lambda_{n,p}); (u_{n,1}, \dots, u_{n,p}))$  denote a system of eigenvalues and eigenvectors of the sample covariance matrix  $S_n \equiv n^{-1}Y_n'Y_n = n^{-1}\Sigma_n^{1/2}X_n'X_n\Sigma_n^{1/2}$ . We can assume that the eigenvalues are sorted in increasing order without loss of generality (w.l.o.g.). The first subscript,  $n$ , will be omitted when no confusion is possible. The e.d.f. of the sample eigenvalues is defined as:  $\forall \lambda \in \mathbb{R}, F_n(\lambda) \equiv p^{-1} \sum_{i=1}^p \mathbb{1}_{[\lambda_i, +\infty)}(\lambda)$ .

In the remainder of the paper, we shall use the notations  $\text{Re}(z)$  and  $\text{Im}(z)$  for the real and imaginary parts, respectively, of a complex number  $z$ , so that

$$\forall z \in \mathbb{C} \quad z = \text{Re}(z) + i \cdot \text{Im}(z) .$$

The Stieltjes transform of a nondecreasing function  $G$  is defined by

$$\forall z \in \mathbb{C}^+ \quad m_G(z) \equiv \int_{-\infty}^{+\infty} \frac{1}{\lambda - z} dG(\lambda) , \quad (2.1)$$

where  $\mathbb{C}^+$  is the half-plane of complex numbers with strictly positive imaginary part. The Stieltjes transform has a well-known inversion formula

$$G(b) - G(a) = \lim_{\eta \rightarrow 0^+} \frac{1}{\pi} \int_a^b \text{Im}[m_G(\xi + i\eta)] d\xi ,$$

which holds if  $G$  is continuous at  $a$  and  $b$ . Thus, the Stieltjes transform of the e.d.f. of sample eigenvalues is

$$\forall z \in \mathbb{C}^+ \quad m_{F_n}(z) = \frac{1}{p} \sum_{i=1}^p \frac{1}{\lambda_i - z} = \frac{1}{p} \text{Tr}[(S_n - zI)^{-1}] ,$$

where  $I$  denotes a conformable identity matrix.

## 2.2 Marčenko-Pastur Equation and Reformulations

Marčenko and Pastur (1967) and others have proven that  $F_n(\lambda)$  converges almost surely (a.s.) to some nonrandom limit  $F(\lambda)$  at all points of continuity of  $F$  under certain sets of assumptions. Furthermore, Marčenko and Pastur discovered the equation that relates  $m_F$  to  $H$ . The most convenient expression of the Marčenko-Pastur equation is the one found in Silverstein (1995, Equation (1.4)):

$$\forall z \in \mathbb{C}^+ \quad m_F(z) = \int_{-\infty}^{+\infty} \frac{1}{\tau[1 - c - cz m_F(z)] - z} dH(\tau) . \quad (2.2)$$

This version of the Marčenko-Pastur equation is the one that we start out with. In addition, Silverstein and Choi (1995) showed that:  $\forall \lambda \in \mathbb{R} - \{0\}$ ,  $\lim_{z \in \mathbb{C}^+ \rightarrow \lambda} m_F(z) \equiv \check{m}_F(\lambda)$  exists, and that  $F$  has a continuous derivative  $F' = \pi^{-1} \text{Im}[\check{m}_F]$  on all of  $\mathbb{R}$  with  $F' \equiv 0$  on  $(-\infty, 0]$ . For purposes that will become apparent later, it is useful to reformulate the Marčenko-Pastur equation.

The limiting e.d.f. of the eigenvalues of  $n^{-1}Y_n'Y_n = n^{-1}\Sigma_n^{1/2}X_n'X_n\Sigma_n^{1/2}$  was defined as  $F$ . In addition, define the limiting e.d.f. of the eigenvalues of  $n^{-1}Y_nY_n' = n^{-1}X_n\Sigma_nX_n'$  as  $\underline{F}$ .

It then holds

$$\begin{aligned}
\forall x \in \mathbb{R} \quad \underline{F}(x) &= (1 - c) \mathbb{1}_{[0, +\infty)}(x) + c F(x) \\
\forall x \in \mathbb{R} \quad F(x) &= \frac{c - 1}{c} \mathbb{1}_{[0, +\infty)}(x) + \frac{1}{c} \underline{F}(x) \\
\forall z \in \mathbb{C}^+ \quad m_{\underline{F}}(z) &= \frac{c - 1}{z} + c m_F(z) \\
\forall z \in \mathbb{C}^+ \quad m_F(z) &= \frac{1 - c}{cz} + \frac{1}{c} m_{\underline{F}}(z) .
\end{aligned}$$

With this notation, Equation (1.3) of Silverstein and Choi (1995) rewrites the Marčenko-Pastur equation as: for each  $z \in \mathbb{C}^+$ ,  $m_{\underline{F}}(z)$  is the unique solution in  $\mathbb{C}^+$  to the equation

$$m_{\underline{F}}(z) = - \left[ z - c \int_{-\infty}^{+\infty} \frac{\tau}{1 + \tau m_{\underline{F}}(z)} dH(\tau) \right]^{-1} . \quad (2.3)$$

Now introduce  $u_{\underline{F}}(z) \equiv -1/m_{\underline{F}}(z)$ . Notice that:  $u_{\underline{F}}(z) \in \mathbb{C}^+ \iff m_{\underline{F}}(z) \in \mathbb{C}^+$ . The mapping from  $u_{\underline{F}}(z)$  to  $m_{\underline{F}}(z)$  is one-to-one on  $\mathbb{C}^+$ .

With this change of variable, Equation (2.3) is equivalent to saying that: for each  $z \in \mathbb{C}^+$ ,  $u_{\underline{F}}(z)$  is the unique solution in  $\mathbb{C}^+$  to the equation

$$u_{\underline{F}}(z) = z + c u_{\underline{F}}(z) \int_{-\infty}^{+\infty} \frac{\tau}{\tau - u_{\underline{F}}(z)} dH(\tau) . \quad (2.4)$$

Let the linear operator  $L$  transform any c.d.f.  $G$  into

$$LG(x) \equiv \int_{-\infty}^x \tau dG(\tau) .$$

Combining  $L$  with the Stieltjes transform, we get

$$m_{LG}(z) = \int_{-\infty}^{+\infty} \frac{\tau}{\tau - z} dG(\tau) = 1 + z m_G(z) .$$

Thus, we can rewrite Equation (2.4) more concisely as

$$u_{\underline{F}}(z) = z + c u_{\underline{F}}(z) m_{LH}(u_{\underline{F}}(z)) . \quad (2.5)$$

As Silverstein and Choi (1995, Equation (1.4)) explain, the function defined in Equation (2.3) is invertible. Thus, we can define the inverse function

$$z_{\underline{F}}(m) \equiv -\frac{1}{m} + c \int_{-\infty}^{+\infty} \frac{\tau}{1 + \tau m} dH(\tau) . \quad (2.6)$$

We can do the same thing for Equation (2.5) and define the inverse function

$$\tilde{z}_{\underline{F}}(u) \equiv u - c u m_{LH}(u) . \quad (2.7)$$

Equations (2.2), (2.3), (2.5), (2.6), and (2.7) are all completely equivalent to one another: solving any one of them means having solved them all. They are all just reformulations of the Marčenko-Pastur equation.

As will be detailed in Section 3, the oracle nonlinear shrinkage estimator of  $\Sigma_n$  involves the quantity  $\check{m}_F(\lambda)$ , for various inputs  $\lambda$ . Subsection 2.3 describes how this quantity can be found in the hypothetical case that  $F$  and  $H$  are actually known. This will then allow us later to discuss consistent estimation of  $\check{m}_F(\lambda)$  in the realistic case when  $F$  and  $H$  are unknown.

### 2.3 Solving the Marčenko-Pastur Equation

Silverstein and Choi (1995) explain how the support of  $F$ , denoted by  $\text{Supp}(F)$  is determined. Let  $B \equiv \{u \in \mathbb{R} : u \neq 0, u \in \text{Supp}^b(H)\}$ . Then plot the function  $\tilde{z}_F(u)$  of (2.7) on the set  $B$ . Find the extreme values on each interval. Delete these points and everything in between on the real line. Do this for all increasing intervals. What is left is just  $\text{Supp}(F)$ ; see Figure 1 of Bai and Silverstein (1998) for an illustration.

To simplify, we will assume from here on that  $\text{Supp}(F)$  is a single compact interval, bounded away from zero, with  $F' > 0$  in the interior of this interval. But if  $\text{Supp}(F)$  is the union of a finite number of such intervals, the arguments presented in this section as well as in the remainder of the paper apply separately to each interval. In particular, our consistency results presented in subsequent sections can be easily extended to this more general case. On the other hand, the even more general case of  $\text{Supp}(F)$  being the union of an infinite number of such intervals or being a non-compact interval is ruled out by Assumption (A4). By our assumption then,  $\text{Supp}(F)$  is given by the compact interval  $[\tilde{z}_F(u_1), \tilde{z}_F(u_2)]$  for some  $u_1 < u_2$ . To keep the notation shorter in what follows, let  $\tilde{z}_1 \equiv \tilde{z}_F(u_1)$  and  $\tilde{z}_2 \equiv \tilde{z}_F(u_2)$ .

We know that for every  $\lambda$  in the interior of  $\text{Supp}(F)$ , there exists a unique  $v \in \mathbb{C}^+$ , denoted by  $v_\lambda$ , such that

$$v_\lambda - c v_\lambda m_{LH}(v_\lambda) = \lambda. \quad (2.8)$$

We further know that

$$F'(\lambda) = \frac{1}{c} \underline{F}'(\lambda) = \frac{1}{c\pi} \text{Im}[\check{m}_F(\lambda)] = \frac{1}{c\pi} \text{Im} \left[ -\frac{1}{v_\lambda} \right].$$

The converse is also true. Since  $\text{Supp}(F) = [\tilde{z}_F(u_1), \tilde{z}_F(u_2)]$ , for every  $x \in (u_1, u_2)$ , there exists a unique  $y > 0$ , denoted by  $y_x$ , such that

$$(x + iy_x) - c(x + iy_x) m_{LH}(x + iy_x) \in \mathbb{R}.$$

In other words,  $y_x$  is the unique value of  $y > 0$  for which  $\text{Im}[(x + iy) - c(x + iy) m_{LH}(x + iy)] = 0$ . Also, if  $\lambda_x$  denotes the value of  $\lambda$  for which we have  $(x + iy_x) - c(x + iy_x) m_{LH}(x + iy_x) = \lambda$ , then, by definition,  $z_{\lambda_x} = x + iy_x$ .

Once we find a way to consistently estimate  $y_x$  for any  $x \in [u_1, u_2]$ , then we have an estimate of the (asymptotic) solution to the Marčenko-Pastur equation. For example,  $\text{Im}[-1/(x + iy_x)]/(c\pi)$  is the value of the density  $F'$  evaluated at  $\text{Re}[(x + iy_x) - c(x + iy_x) m_{LH}(x + iy_x)] = (x + iy_x) - c(x + iy_x) m_{LH}(x + iy_x)$ .

From the above arguments, it follows that

$$\forall \lambda \in (\tilde{z}_1, \tilde{z}_2) \quad \check{m}_F(\lambda) = -\frac{1}{v_\lambda} \quad \text{and so} \quad \check{m}_F(\lambda) = \frac{1-c}{c\lambda} - \frac{1}{c} \frac{1}{v_\lambda}. \quad (2.9)$$



### 3 Oracle Estimator

### 3.1 Covariance Matrix

In the absence of specific information about the true covariance matrix  $\Sigma_n$ , it appears reasonable to restrict attention to the class of estimators that are equivariant with respect to rotations of the observed data. To be more specific, let  $W$  be an arbitrary  $p$ -dimensional orthogonal matrix. Let  $\hat{\Sigma}_n \equiv \hat{\Sigma}_n(Y_n)$  be an estimator of  $\Sigma_n$ . Then the estimator is said to be *rotation-equivariant* if it satisfies  $\hat{\Sigma}_n(Y_n W) = W' \hat{\Sigma}_n(Y_n) W$ . In other words, the estimate based on the rotated data equals the rotation of the estimate based on the original data. The class of rotation-equivariant estimators of the covariance matrix is constituted of all the estimators that have the same eigenvectors as the sample covariance matrix; for example, see Perlman (2007, Section 5.4). Every rotation-equivariant estimator is thus of the form

$U_n D_n U'_n$  where  $D_n \equiv \text{Diag}(d_1, \dots, d_p)$  is diagonal,

and where  $U_n$  is the matrix whose  $i^{\text{th}}$  column is the sample eigenvector  $u_i \equiv u_{n,i}$ . This is the class we consider.

The starting objective is to find the matrix in this class that is closest to  $\Sigma_n$ . To measure distance, we choose the Frobenius norm defined as

$$\|A\| \equiv \sqrt{\text{Tr}(AA')}/r \quad \text{for any matrix } A \text{ of dimension } r \times m. \quad (3.1)$$

(Dividing by the dimension of the square matrix  $AA'$  inside the root is not standard, but we do this for asymptotic purposes so that the Frobenius norm remains constant equal to one for the identity matrix regardless of the dimension; see Ledoit and Wolf (2004).) As a result, we end up with the following minimization problem:

$$\min_{D_n} ||U_n D_n U_n' - \Sigma_n|| \text{ .}$$

Elementary matrix algebra shows that its solution is

$$D_n^* \equiv \text{Diag}(d_1^*, \dots, d_p^*) \quad \text{where} \quad d_i^* \equiv u_i' \Sigma_n u_i \quad \text{for } i = 1, \dots, p. \quad (3.2)$$

The interpretation of  $d_i^*$  is that it captures how the  $i^{\text{th}}$  sample eigenvector  $u_i$  relates to the population covariance matrix  $\Sigma_n$  as a whole. As a result, the finite-sample optimal estimator is given by

$$S_n^* \equiv U_n D_n^* U_n' \quad \text{where} \quad D_n^* \text{ is defined as in (3.2) .} \quad (3.3)$$

By generalizing the Marčenko-Pastur equation (2.2), Ledoit and Pécché (2011) show that  $d_i^*$  can be approximated by the quantity

$$d_i^{or} \equiv \frac{\lambda_i}{|1 - c - c \lambda_i \check{m}_F(\lambda_i)|^2} \quad \text{for } i = 1, \dots, p, \quad (3.4)$$

from which they deduce their oracle estimator

$$S_n^{or} \equiv U_n D_n^{or} U_n' \quad \text{where} \quad D_n^{or} \equiv \text{Diag}(d_1^{or}, \dots, d_p^{or}). \quad (3.5)$$

The key difference between  $D_n^*$  and  $D_n^{or}$  is that the former depends on the unobservable population covariance matrix, whereas the latter depends on the limiting distribution of sample eigenvalues, which makes it amenable to estimation, as explained below.

Note that  $S_n^{or}$  constitutes a nonlinear shrinkage estimator: since the value of the denominator of  $d_i^{or}$  varies with  $\lambda_i$ , the shrunk eigenvalues  $d_i^{or}$  are obtained by applying a nonlinear transformation to the sample eigenvalues  $\lambda_i$ ; see Figure 3 for an illustration. Ledoit and P  ch   (2011) also illustrate in some (limited) simulations that this oracle estimator can provide a magnitude of improvement over the linear shrinkage estimator of Ledoit and Wolf (2004).

### 3.2 Precision Matrix

Often times an estimator of the inverse of the covariance matrix, or the precision matrix,  $\Sigma_n^{-1}$  is required. A reasonable strategy would be to first estimate  $\Sigma_n$  and to then simply take the inverse of the resulting estimator. However, such a strategy will generally not be optimal.

By arguments analogous to those leading up to (3.3), among the class of rotation-equivariant estimators, the finite-sample optimal estimator of  $\Sigma_n^{-1}$  with respect to the Frobenius norm is given by

$$P_n^* \equiv U_n A_n^* U_n' \quad \text{where} \quad a_i^* \equiv u_i' \Sigma_n^{-1} u_i \quad \text{for } i = 1, \dots, p. \quad (3.6)$$

In particular, note that  $P_n^* \neq (S_n^*)^{-1}$  in general.

Studying the asymptotic behavior of the diagonal matrix  $A_n^*$  led Ledoit and P  ch   (2011) to the following oracle estimator:

$$P_n^{or} \equiv U_n A_n^{or} U_n' \quad \text{where} \quad a_i^{or} \equiv \lambda_i^{-1} (1 - c - 2c\lambda_i \operatorname{Re}[\check{m}_F(\lambda_i)]) \quad \text{for } i = 1, \dots, p. \quad (3.7)$$

In particular, note that  $P_n^{or} \neq (S_n^{or})^{-1}$  in general.

**Remark 3.1.** One can see that both oracle estimators  $S_n^{or}$  and  $P_n^{or}$  involve the unknown quantities  $\check{m}_F(\lambda_i)$ , for  $i = 1, \dots, p$ . As a result, they are not *bona fide* estimators. However, being able to consistently estimate  $\check{m}_F(\lambda)$ , uniformly in  $\lambda$ , will allow us to construct *bona fide* estimators  $\hat{S}_n$  and  $\hat{P}_n$  that converge to their respective oracle counterparts almost surely (in the sense that the Frobenius norm of the difference converges to zero almost surely).

Section 4 explains how to construct a uniformly consistent estimator of  $\check{m}_F(\lambda)$  based on a consistent estimator of  $H$ , the limiting spectral distribution of the population eigenvalues. Section 5 discusses how to construct a consistent estimator of  $H$  from the data. ■

### 3.3 Further Details on the Results of Ledoit and P  ch   (2011)

Ledoit and P  ch   (2011) (hereafter LP) study functionals of the type

$$\forall z \in \mathbb{C}^+, \quad \Theta_N^g(z) \equiv \frac{1}{N} \sum_{i=1}^N \frac{1}{\lambda_i - z} \sum_{j=1}^N |u_i^* v_j|^2 \times g(\tau_j) = \frac{1}{N} \operatorname{Tr} [(S_N - zI)^{-1} g(\Sigma_N)], \quad (3.8)$$

where  $g$  is any real-valued univariate function satisfying suitable regularity conditions. Comparison with Equation (2.1) reveals that this family of functionals generalizes the Stieltjes

transform, with the Stieltjes transform corresponding to the special case  $g \equiv 1$ . What is of interest is what happens for other, non-constant functions  $g$ .

It turns out that it is possible to generalize the Marčenko-Pastur result (2.2) to any function  $g$  with finitely many points of discontinuity. Under assumptions that are usual in the Random Matrix Theory literature, LP prove in their Theorem 2 that there exists a non-random function  $\Theta^g$  defined over  $\mathbb{C}^+$  such that  $\Theta_N^g(z)$  converges a.s. to  $\Theta^g(z)$  for all  $z \in \mathbb{C}^+$ . Furthermore,  $\Theta^g$  is given by

$$\forall z \in \mathbb{C}^+ \quad \Theta^g(z) \equiv \int_{-\infty}^{+\infty} \frac{g(\tau)}{\tau[1 - c - cz m_F(z)] - z} dH(\tau) . \quad (3.9)$$

What is remarkable is that, as one moves from the constant function  $g \equiv 1$  to any other function  $g(\tau)$ , the integration kernel  $\frac{g(\tau)}{\tau[1 - c - cz m_F(z)] - z}$  remains unchanged. Therefore Equation (3.9) is a direct generalization of Marčenko and Pastur's foundational result.

The power and usefulness of this generalization become apparent once one starts plugging specific, judiciously chosen functions  $g(\tau)$  into Equation (3.9). For the purpose of illustration, LP work out three examples of functions  $g(\tau)$ .

The first example of LP is  $g(\tau) \equiv \mathbb{1}_{(-\infty, \tau)}$ , where  $\mathbb{1}$  denotes the indicator function of a set. It enables them to characterize the asymptotic location of sample eigenvectors relative to population eigenvectors. Since this result is not directly relevant to the present paper, we will not elaborate further, and refer the interested reader to LP's Section 1.2.

The second example of LP is  $g(\tau) \equiv \tau$ . It enables them to characterize the asymptotic behavior of the quantities  $d_i^{or}$  introduced in Equation (3.4). More formally, for any  $u \in (0, 1)$  define

$$\Delta_n^*(u) \equiv \frac{1}{p} \sum_{i=1}^{\lfloor u \cdot p \rfloor} d_i^* \quad \text{and} \quad \Delta_n^{or}(u) \equiv \frac{1}{p} \sum_{i=1}^{\lfloor u \cdot p \rfloor} d_i^{or} , \quad (3.10)$$

where  $\lfloor \cdot \rfloor$  denotes the integer part. LP's Theorem 4 proves that  $\Delta_n^*(u) - \Delta_n^{or}(u) \rightarrow 0$  a.s.

The third example of LP is  $g(\tau) \equiv 1/\tau$ . It enables them to characterize the asymptotic behavior of the quantities  $a_i^{or}$  introduced in Equation (3.7). For any  $u \in (0, 1)$  define

$$\Psi_n^*(u) \equiv \frac{1}{p} \sum_{i=1}^{\lfloor u \cdot p \rfloor} a_i^* \quad \text{and} \quad \Psi_n^{or}(u) \equiv \frac{1}{p} \sum_{i=1}^{\lfloor u \cdot p \rfloor} a_i^{or} . \quad (3.11)$$

LP's Theorem 5 proves that  $\Psi_n^*(u) - \Psi_n^{or}(u) \rightarrow 0$  a.s.

## 4 Estimation of $\check{m}_F(\lambda)$

Fix  $x \in [u_1 + \eta, u_2 - \eta]$ , where  $\eta > 0$  is some small number. From the previous discussion in Section 2, it follows that the equation

$$\text{Im}[x + iy - c(x + iy) m_{LH}(x + iy)] = 0$$

has a unique solution  $y \in (0, +\infty)$ , called  $y_x$ . Since  $u_1 < x < u_2$ , it follows that  $y_x > 0$ ; for  $x = u_1$  or  $x = u_2$ , we would have  $y_x = 0$  instead. The goal is to consistently estimate  $y_x$ , uniformly in  $x \in [u_1 + \eta, u_2 - \eta]$ .

Define for any c.d.f.  $G$  and for any  $d > 0$ , the real function

$$g_{G,d}(y, x) \equiv \left| \operatorname{Im} [x + iy - d(x + iy) m_{LG}(x + iy)] \right|.$$

With this notation,  $y_x$  is the unique minimizer in  $(0, +\infty)$  of  $g_{H,c}(y, x)$  then. In particular,  $g_{H,c}(y_x, x) = 0$ .

In the remainder of the paper, the symbol  $\Rightarrow$  denotes weak convergence (or convergence in distribution).

**Proposition 4.1.**

(i) Let  $\{\hat{H}_n\}$  be a sequence of probability measures with  $\hat{H}_n \Rightarrow H$ . Let  $\{\hat{c}_n\}$  be a sequence of positive real numbers with  $\hat{c}_n \rightarrow c$ . Let  $K \subseteq (0, \infty)$  be a compact interval satisfying  $\{y_x : x \in [u_1 + \eta, u_2 - \eta]\} \subseteq K$ . For a given  $x \in [u_1 + \eta, u_2 - \eta]$ , let  $\hat{y}_{n,x} \equiv \min_{y \in K} g_{\hat{H}_n, \hat{c}_n}(y, x)$ . It then holds that  $\hat{y}_{n,x} \rightarrow y_x$  uniformly in  $x \in [u_1 + \eta, u_2 - \eta]$ .

(ii) In case of  $\hat{H}_n \Rightarrow H$  a.s., it holds that  $\hat{y}_{n,x} \rightarrow y_x$  a.s. uniformly in  $x \in [u_1 + \eta, u_2 - \eta]$ .

It should be pointed out that the assumption  $\{y_x : x \in [u_1 + \eta, u_2 - \eta]\} \subseteq K$  is not really restrictive, since one can choose  $K \equiv [\varepsilon, 1/\varepsilon]$ , for  $\varepsilon$  arbitrarily small.

We also need to solve the ‘inverse’ estimation problem, namely starting with  $\lambda$  and recovering the corresponding  $v_\lambda$ . Fix  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ , where  $\tilde{\delta} > 0$  is some small number. From the previous discussion, it follows that the equation

$$v - c v m_{LH}(v) = \lambda$$

has a unique solution  $v \in \mathbb{C}^+$ , called  $v_\lambda$ . The goal is to consistently estimate  $v_\lambda$ , uniformly in  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ .

Define for any c.d.f.  $G$  and for any  $d > 0$ , the real function

$$h_{G,d}(v, \lambda) \equiv |v - d v m_{LG}(v) - \lambda|.$$

With this notation,  $v_\lambda$  is the unique minimizer in  $\mathbb{C}^+$  of  $h_{H,c}(v, \lambda)$  then. In particular,  $h_{H,c}(v_\lambda, \lambda) = 0$ .

**Proposition 4.2.**

(i) Let  $\{\hat{H}_n\}$  be a sequence of probability measures with  $\hat{H}_n \Rightarrow H$ . Let  $\{\hat{c}_n\}$  be a sequence of positive real numbers with  $\hat{c}_n \rightarrow c$ . Let  $K \subseteq \mathbb{C}^+$  be a compact set satisfying  $\{v_\lambda : \lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]\} \subseteq K$ . For a given  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ , let  $\hat{v}_{n,\lambda} \equiv \min_{v \in K} h_{\hat{H}_n, \hat{c}_n}(v, \lambda)$ . It then holds that  $\hat{v}_{n,\lambda} \rightarrow v_\lambda$  uniformly in  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ .

(ii) In case of  $\hat{H}_n \Rightarrow H$  a.s., it holds that  $\hat{v}_{n,\lambda} \rightarrow v_\lambda$  a.s. uniformly in  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ .

Being able to find consistent estimators of  $v_\lambda$ , uniformly in  $\lambda$ , now allows us to find consistent estimators of  $\check{m}_F(\lambda)$ , uniformly in  $\lambda$ , based on (2.9). Our estimator of  $\check{m}_F(\lambda)$  is given by

$$\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda) \equiv \frac{1 - \hat{c}_n}{\hat{c}_n \lambda} - \frac{1}{\hat{c}_n} \frac{1}{\hat{v}_{n, \lambda}}. \quad (4.1)$$

This, in turn, provides us with a consistent estimator of  $S_n^{or}$ , the oracle nonlinear shrinkage estimator of  $\Sigma_n$ . Define

$$\hat{S}_n \equiv U_n \hat{D}_n U_n' \quad \text{where} \quad \hat{d}_i \equiv \frac{\lambda_i}{|1 - \hat{c}_n - \hat{c}_n \lambda_i \check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda_i)|^2} \quad \text{for } i = 1, \dots, p. \quad (4.2)$$

It also provides us with a consistent estimator of  $P_n^{or}$ , the oracle nonlinear shrinkage estimator of  $\Sigma_n^{-1}$ . Define

$$\hat{P}_n \equiv U_n \hat{A}_n U_n' \quad \text{where} \quad \hat{a}_i \equiv \lambda_i^{-1} (1 - \hat{c}_n - 2 \hat{c}_n \lambda_i \operatorname{Re}[\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda_i)]) \quad \text{for } i = 1, \dots, p. \quad (4.3)$$

In particular, note that  $\hat{P}_n \neq \hat{S}_n^{-1}$  in general.

**Proposition 4.3.**

(i) Let  $\{\hat{H}_n\}$  be a sequence of probability measures with  $\hat{H}_n \Rightarrow H$ . Let  $\{\hat{c}_n\}$  be a sequence of positive real numbers with  $\hat{c}_n \rightarrow c$ . It then holds that:

- (a)  $\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda) \rightarrow \check{m}_F(\lambda)$  uniformly in  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$
- (b)  $\|\hat{S}_n - S_n^{or}\| \rightarrow 0$
- (c)  $\|\hat{P}_n - P_n^{or}\| \rightarrow 0$

(ii) In case of  $\hat{H}_n \Rightarrow H$  a.s., it holds that:

- (a)  $\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda) \rightarrow \check{m}_F(\lambda)$  uniformly in  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$  a.s.
- (b)  $\|\hat{S}_n - S_n^{or}\| \rightarrow 0$  a.s.
- (c)  $\|\hat{P}_n - P_n^{or}\| \rightarrow 0$  a.s.

## 5 Estimation of $H$

As described before, consistent estimation of the oracle estimators of Ledoit and P  ch   (2011) requires (uniformly) consistent estimation of  $\check{m}_F(\lambda)$ . Since  $\operatorname{Im}[\check{m}_F(\lambda)] = \pi F'(\lambda)$ , one possible approach could be to take an off-the-shelf density estimator for  $F'$ , based on the observed sample eigenvalues  $\lambda_i$ . There exists a large literature on density estimation; for example, see Silverman (1986). The real part of  $\check{m}_F(\lambda_i)$  could be estimated in a similar manner.

However, the sample eigenvalues do not satisfy any of the regularity conditions usually invoked for the underlying data. It really is not clear at all whether an off-the-shelf density estimator applied to the sample eigenvalues would result in consistent estimation of  $F'$ .

Even if this issue was somehow resolved, using such a generic procedure would not exploit the specific features of the problem. Namely:  $F$  is not just any distribution, it is a distribution of sample eigenvalues. It is the solution to the Marčenko-Pastur equation for some  $H$ . This is valuable information that narrows down considerably the set of possible distributions  $F$ . Therefore an estimation procedure specifically designed to incorporate this *a priori* knowledge would be better suited to the problem at hand. This is the approach we select.

In a nutshell: our estimator of  $F$  is the c.d.f. that is closest to  $F_n$  among the c.d.f.s that are a solution to the Marčenko-Pastur equation for some  $\tilde{H}$  and for  $\tilde{c} \equiv \hat{c}_n \equiv p/n$ . The ‘underlying’ distribution  $\tilde{H}$  which produces the thus obtained estimator of  $F$  is, in turn, our estimator of  $H$ . If we can show that this estimator of  $H$  is consistent, then the results of the previous section demonstrate that the implied estimator of  $\check{m}_F(\lambda)$  is uniformly consistent.

Subsection 5.1 derives theoretical properties of this approach, while Subsection 5.2 discusses various issues concerning the practical implementation.

## 5.1 Consistency Results

For a grid of real numbers  $Q \equiv \{\dots, t_{-1}, t_0, t_1, \dots\} \subseteq \mathbb{R}$ , with  $t_{k-1} < t_k$ , define the corresponding grid size  $\gamma$  as

$$\gamma \equiv \sup_k (t_k - t_{k-1}) .$$

A grid  $Q$  is said to cover a compact interval  $[a, b] \subseteq \mathbb{R}$  if there exists at least one  $t_k \in Q$  with  $t_k \leq a$  and at least another  $t_{k'} \in Q$  with  $b \leq t_{k'}$ . A sequence of grids  $\{Q_n\}$  is said to eventually cover a compact interval  $[a, b]$  if for every  $\phi > 0$  there exist  $N \equiv N(\phi)$  such that  $Q_n$  covers the compact interval  $[a + \phi, b - \phi]$  for all  $n \geq N$ .

For any probability measure  $\tilde{H}$  on the real line and for any  $\tilde{c} > 0$ , let  $F_{\tilde{H}, \tilde{c}}$  denote the c.d.f. on the real line induced by the corresponding solution of the Marčenko-Pastur equation. More specifically, for each  $z \in \mathbb{C}^+$ ,  $m_{F_{\tilde{H}, \tilde{c}}}(z)$  is the unique solution for  $m \in \mathbb{C}^+$  to the equation

$$m = \int_{-\infty}^{+\infty} \frac{1}{\tau [1 - \tilde{c} - \tilde{c} z m] - z} d\tilde{H}(\tau) .$$

In this notation, we then have  $F = F_{H, c}$ .

It follows from Silverstein and Choi (1995) again that  $\forall \lambda \in \mathbb{R} - \{0\}$ ,  $\lim_{z \in \mathbb{C}^+ \rightarrow \lambda} m_{F_{\tilde{H}, \tilde{c}}}(z) \equiv \check{m}_{F_{\tilde{H}, \tilde{c}}}(\lambda)$  exists, and that  $F_{\tilde{H}, \tilde{c}}$  has a continuous derivative  $F'_{\tilde{H}, \tilde{c}} = \pi^{-1} \text{Im} \left[ \check{m}_{F_{\tilde{H}, \tilde{c}}} \right]$  on  $(0, +\infty)$ . In the case  $\tilde{c} < 1$ ,  $F_{\tilde{H}, \tilde{c}}$  has a continuous derivative on all of  $\mathbb{R}$  with  $F'_{\tilde{H}, \tilde{c}} \equiv 0$  on  $(-\infty, 0]$ .

For a grid  $Q$  on the real line and for two c.d.f.s  $G_1$  and  $G_2$ , define

$$\|G_1 - G_2\|_Q \equiv \sup_{t \in Q} |G_1(t) - G_2(t)| .$$

The following theorem shows that both  $F$  and  $H$  can be estimated consistently via an idealized algorithm.

**Theorem 5.1.** *Let  $\{Q_n\}$  be a sequence of grids on the real line eventually covering the support of  $F$  with corresponding grid sizes  $\{\gamma_n\}$  satisfying  $\gamma_n \rightarrow 0$ . Let  $\{\widehat{c}_n\}$  be a sequence of positive real numbers with  $\widehat{c}_n \rightarrow c$ . Let  $\widehat{H}_n$  be defined as*

$$\widehat{H}_n \equiv \operatorname{argmin}_{\widetilde{H}} \|F_{\widetilde{H}, \widehat{c}_n} - F_n\|_{Q_n}, \quad (5.1)$$

where  $\widetilde{H}$  is a probability measure.

Then we have (i)  $F_{\widehat{H}_n, \widehat{c}_n} \Rightarrow F$  a.s.; and (ii)  $\widehat{H}_n \Rightarrow H$  a.s.

The algorithm used in the theorem is not practical for two reasons. First, it is not possible to optimize over all probability measures  $\widetilde{H}$ . But similarly to El Karoui (2008), we can show that it is sufficient to optimize over all probability measures which are sums of atoms, the location of which is restricted to a fixed-size grid, with the grid size vanishing asymptotically.

**Corollary 5.1.** *Let  $\{Q_n\}$  be a sequence of grids on the real line eventually covering the support of  $F$  with corresponding grid sizes  $\{\gamma_n\}$  satisfying  $\gamma_n \rightarrow 0$ . Let  $\{\widehat{c}_n\}$  be a sequence of positive real numbers with  $\widehat{c}_n \rightarrow c$ . Let  $\mathcal{P}_n$  denote the set of all probability measures which are sums of atoms belonging to the grid  $\{J_n/T_n, (J_n+1)/T_n, \dots, K_n/T_n\}$  with  $T_n \rightarrow \infty$ ,  $J_n$  being the largest integer satisfying  $J_n/T_n \leq \lambda_1$ , and  $K_n$  being the smallest integer satisfying  $K_n/T_n \geq \lambda_p$ . Let  $\widehat{H}_n$  be defined as*

$$\widehat{H}_n \equiv \operatorname{argmin}_{\widetilde{H} \in \mathcal{P}_n} \|F_{\widetilde{H}, \widehat{c}_n} - F_n\|_{Q_n}, \quad (5.2)$$

Then we have (i)  $F_{\widehat{H}_n, \widehat{c}_n} \Rightarrow F$  a.s.; and (ii)  $\widehat{H}_n \Rightarrow H$  a.s.

But even restricting the optimization over a manageable set of probability measures is not quite practical yet for a second reason. Namely, to compute  $F_{\widetilde{H}, \widehat{c}_n}$  exactly for a given  $\widetilde{H}$ , one would have to (numerically) solve the Marčenko-Pastur equation for an infinite number of points. In practice, we can only afford to solve the equation for a finite number of points and then approximate  $F_{\widetilde{H}, \widehat{c}_n}$  by trapezoidal integration. Fortunately, this approximation does not negatively affect the consistency of our estimators.

Let  $G$  be a c.d.f. with continuous density  $g$  and compact support  $[a, b]$ . For a grid  $Q \equiv \{\dots, t_{-1}, t_0, t_1, \dots\}$  covering the support of  $G$ , the approximation to  $G$  via trapezoidal integration over the grid  $Q$ , denoted by  $\widehat{G}_Q$ , is obtained as follows. For  $t \in [a, b]$ , let  $J_{lo} \equiv \max\{k : t_k \leq a\}$  and  $J_{hi} \equiv \min\{k : t < t_k\}$ . Then

$$\widehat{G}_Q(t) \equiv \sum_{k=J_{lo}}^{J_{hi}-1} \frac{(t_{k+1} - t_k)[g(t_k) + g(t_{k+1})]}{2}. \quad (5.3)$$

Now turn to the special case  $G \equiv F_{\widetilde{H}, \widehat{c}}$  and  $Q \equiv Q_n$ . In this case, we denote the approximation to  $F_{\widetilde{H}, \widehat{c}}$  via trapezoidal integration over the grid  $Q_n$  by  $\widehat{F}_{\widetilde{H}, \widehat{c}; Q_n}$ .

**Corollary 5.2.** *Assume the same assumptions as in Corollary 5.1. Let  $\widehat{H}_n$  be defined as*

$$\widehat{H}_n \equiv \operatorname{argmin}_{\widetilde{H} \in \mathcal{P}_n} \|\widehat{F}_{\widetilde{H}, \widehat{c}; Q_n} - F_n\|_{Q_n}, \quad (5.4)$$

Let  $\check{m}_{F_{\widehat{H}_n, \widehat{c}_n}}(\lambda)$ ,  $\widehat{S}_n$ , and  $\widehat{P}_n$  be defined as in (4.1), (4.2), and (4.3), respectively. Then:

- (i)  $F_{\hat{H}_n, \hat{c}_n} \Rightarrow F$  a.s.
- (ii)  $\hat{H}_n \Rightarrow H$  a.s.
- (iii) For any  $\tilde{\delta} > 0$ ,  $\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda) \rightarrow \check{m}_F(\lambda)$  a.s. uniformly in  $\lambda \in [\tilde{z}_1 - \tilde{\delta}, \tilde{z}_2 + \tilde{\delta}]$ .
- (iv)  $\|\hat{S}_n - S_n^{or}\| \rightarrow 0$  a.s.
- (v)  $\|\hat{P}_n - P_n^{or}\| \rightarrow 0$  a.s.

## 5.2 Implementation Details

**Decomposition of the c.d.f. of Population Eigenvalues** As discussed before, it is not practical to search over the set of all possible c.d.f.s  $\tilde{H}$ . Following El Karoui (2008), we project  $H$  onto a certain basis of c.d.f.s  $(M_k)_{k=1, \dots, K}$ , where  $K$  goes to infinity along with  $n$  and  $p$ . The projection of  $H$  onto this basis is given by the nonnegative weights  $w_1, \dots, w_K$ , where

$$\forall t \in \mathbb{R} \quad H(t) \approx \tilde{H}(t) \equiv \sum_{k=1}^K w_k M_k(t) \quad \text{and} \quad \sum_{k=1}^K w_k = 1. \quad (5.5)$$

Thus, our estimator for  $F$  will be a solution to the Marčenko-Pastur equation for  $\tilde{H}$  given by Equation (5.5) for some  $(w_k)_{k=1, \dots, K}$ , and for  $\tilde{c} \equiv p/n$ . It is just a matter of searching over all sets of nonnegative weights summing up to one.

**Choice of Basis** We base the c.d.f.s  $(M_k)_{k=1, \dots, K}$  on a grid of  $p$  equally spaced points on the interval  $[\lambda_1, \lambda_p]$ :

$$x_i \equiv \lambda_1 + \frac{i-1}{p}(\lambda_p - \lambda_1) \quad \text{for } i = 1, \dots, p. \quad (5.6)$$

Thus,  $x_1 = \lambda_1$  and  $x_p = \lambda_p$ . We then form the basis  $\{M_1, \dots, M_K\}$  as the union of three families of c.d.f.s:

1. the indicator functions  $\mathbb{1}_{[x_i, +\infty)}$  ( $i = 1, \dots, p$ );
2. the c.d.f.s whose derivatives are linearly increasing on the interval  $[x_{i-1}, x_i]$  and zero everywhere else ( $i = 2, \dots, p$ );
3. the c.d.f.s whose derivatives are linearly decreasing on the interval  $[x_{i-1}, x_i]$  and zero everywhere else ( $i = 2, \dots, p$ ).

This list yields a basis  $(M_k)_{k=1, \dots, K}$  of dimension  $K = 3p - 2$ . Notice that by the theoretical results of Section 5.1, it would be sufficient to use the first family only. Including the second and third families in addition cannot make the approximation to  $H$  any worse.



**Trapezoidal Integration** For a given  $\tilde{H} \equiv \sum_{k=1}^K w_k M_k$ , it is computationally too expensive (in the context of an optimization procedure) to solve the Marčenko-Pastur equation for  $m_F(z)$  over all  $z \in \mathbb{C}^+$ . It is more efficient to solve the Marčenko-Pastur equation only for  $\check{m}_F(x_i)$  ( $i = 1, \dots, p$ ), and to use the trapezoidal approximation formula to deduce from it  $F(x_i)$  ( $i = 1, \dots, p$ ). The trapezoidal rule gives

$$\begin{aligned} \forall i = 1, \dots, p \quad F(x_i) &= \sum_{j=1}^{i-1} \frac{x_{j+1} - x_{j-1}}{2} F'(x_j) + \frac{x_i - x_{i-1}}{2} F'(x_i) \\ &= \sum_{j=1}^{i-1} \frac{(x_{j+1} - x_{j-1}) \operatorname{Im}[\check{m}_F(x_j)]}{2\pi} + \frac{(x_i - x_{i-1}) \operatorname{Im}[\check{m}_F(x_i)]}{2\pi}, \end{aligned} \quad (5.7)$$

with the convention  $x_0 \equiv 0$ .

**Objective Function** The objective function measures the distance between  $F_n$  and the  $F$  that solves the Marčenko-Pastur equation for  $\tilde{H} \equiv \sum_{k=1}^K w_k M_k$  and for  $\tilde{c} \equiv p/n$ . Traditionally,  $F_n$  is defined as càdlàg, that is :  $F_n(\lambda_1) = 1/p$  and  $F_n(\lambda_p) = 1$ . However, there is a certain degree of arbitrariness in this convention: why is  $F_n(\lambda_p)$  equal to one but  $F_n(\lambda_1)$  not equal to zero? By symmetry, there is no *a priori* justification for specifying that the largest eigenvalue is closer to the supremum of the support of  $F$  than the smallest to its infimum. Therefore, a different convention might be more appropriate in this case, which leads us to the following definition:

$$\forall i = 1, \dots, p \quad \hat{F}_n(\lambda_i) \equiv \frac{i}{p} - \frac{1}{2p}. \quad (5.8)$$

This choice restores a certain element of symmetry to the treatment of the smallest vs. the largest eigenvalue. From Equation (5.8), we deduce  $\hat{F}_n(x_i)$ , for  $i = 2, \dots, p-1$ , by linear interpolation. With a sup-norm error penalty, this leads to the following objective function:

$$\max_{i=1, \dots, p} \left| F(x_i) - \hat{F}_n(x_i) \right|, \quad (5.9)$$

where  $F(x_i)$  is given by Equation (5.7) for  $i = 1, \dots, p$ . Using Equation (5.7), we can rewrite this objective function as

$$\max_{i=1, \dots, p} \left| \sum_{j=1}^{i-1} \frac{(x_{j+1} - x_{j-1}) \operatorname{Im}[\check{m}_F(x_j)]}{2\pi} + \frac{(x_i - x_{i-1}) \operatorname{Im}[\check{m}_F(x_i)]}{2\pi} - \hat{F}_n(x_i) \right|.$$

**Optimization Program** We now have all the ingredients needed to state the optimization program that will extract the estimator of  $\check{m}_F(x_1), \dots, \check{m}_F(x_p)$  from the observations  $\lambda_1, \dots, \lambda_p$ . It is the following:

$$\min_{\substack{m_1, \dots, m_p \\ w_1, \dots, w_K}} \max_{i=1, \dots, p} \left| \sum_{j=1}^{i-1} \frac{(x_{j+1} - x_{j-1}) \operatorname{Im}[m_j]}{2\pi} + \frac{(x_i - x_{i-1}) \operatorname{Im}[m_i]}{2\pi} - \hat{F}_n(x_i) \right|$$

subject to

$$\begin{aligned}
\forall j = 1, \dots, p \quad m_j &= \sum_{k=1}^K \int_{-\infty}^{+\infty} \frac{w_k}{t [1 - (p/n) - (p/n) x_j m_j] - x_j} dM_k(t) \\
\sum_{k=1}^K w_k &= 1 \\
\forall j = 1, \dots, p \quad m_j &\in \mathbb{C}^+ \\
\forall k = 1, \dots, K \quad w_k &\geq 0.
\end{aligned} \tag{5.10}$$

The key is to introduce the variables  $m_j \equiv \check{m}_F(x_j)$ , for  $j = 1, \dots, p$ . The constraint in Equation (5.10) imposes that  $m_j$  is the solution to the Marčenko-Pastur equation evaluated as  $z \in \mathbb{C}^+ \rightarrow x_j$  when  $\tilde{H} = \sum_{k=1}^K w_k M_k$ .

**Real Optimization Program** In practice, most optimizers only accept real variables, therefore it is necessary to decompose  $m_j$  into its real and imaginary parts:  $a_j \equiv \text{Re}[m_j]$  and  $b_j \equiv \text{Im}[m_j]$ . Then we can optimize separately over the two sets of real variables  $a_j$  and  $b_j$  for  $j = 1, \dots, p$ . The Marčenko-Pastur constraint in Equation (5.10) splits into two constraints: one for the real part and the other for the imaginary part. The reformulated optimization program is

$$\min_{\substack{a_1, \dots, a_p \\ b_1, \dots, b_p \\ w_1, \dots, w_K}} \max_{i=1, \dots, p} \left| \sum_{j=1}^{i-1} \frac{(x_{j+1} - x_{j-1}) b_j}{2\pi} + \frac{(x_i - x_{i-1}) b_i}{2\pi} - \hat{F}_n(x_i) \right| \tag{5.11}$$

subject to

$$\forall j = 1, \dots, p \quad a_j = \sum_{k=1}^K \int_{-\infty}^{+\infty} \text{Re} \left\{ \frac{w_k}{t [1 - (p/n) - (p/n) x_j (a_j + ib_j)] - x_j} \right\} dM_k(t) \tag{5.12}$$

$$\forall j = 1, \dots, p \quad b_j = \sum_{k=1}^K \int_{-\infty}^{+\infty} \text{Im} \left\{ \frac{w_k}{t [1 - (p/n) - (p/n) x_j (a_j + ib_j)] - x_j} \right\} dM_k(t) \tag{5.13}$$

$$\sum_{k=1}^K w_k = 1 \tag{5.14}$$

$$\forall j = 1, \dots, p \quad b_j \geq 0 \tag{5.15}$$

$$\forall k = 1, \dots, K \quad w_k \geq 0. \tag{5.16}$$

**Remark 5.1.** Since the theory of Sections 4 and 5.1 partly assumes that  $m_j$  belongs to a compact set in  $\mathbb{C}^+$  bounded away from the real line, we might want to add to the real optimization program the constraints that  $-1/\varepsilon \leq a_j \leq 1/\varepsilon$  and that  $\varepsilon \leq b_j \leq 1/\varepsilon$ , for some small  $\varepsilon > 0$ . Our simulations indicate that for a small value of  $\varepsilon$  such as  $\varepsilon = 10^{-6}$ , this makes no difference in practice.

**Sequential Linear Programming** While the optimization program defined in Equations (5.11)–(5.16) may appear daunting at first sight because of its non-convexity, it is in fact solved quickly and efficiently by off-the-shelf optimization software implementing Sequential Linear Programming (SLP). The key is to linearize Equations (5.12)–(5.13), the two constraints that embody the Marčenko-Pastur equation, around an approximate solution point. Once they are linearized, the optimization program (5.11)–(5.16) becomes a standard Linear Programming (LP) problem, which can be solved very quickly. Then we linearize again Equations (5.12)–(5.13) around the new point, and this generates a new LP problem; hence the name: *Sequential* Linear Programming. The software iterates until a satisfactory degree of convergence is achieved. All of this is handled automatically by the SLP optimizer. The user only needs to specify the problem (5.11)–(5.16), as well as some starting point, and then launch the SLP optimizer. For our SLP optimizer, we selected a standard off-the-shelf commercial software: SNOPT<sup>TM</sup> Version 7.2-5; see Gill et al. (2002). While SNOPT<sup>TM</sup> was originally designed for Sequential Quadratic Programming, it also handles SLP, since Linear Programming can be viewed as a particular case of Quadratic Programming with no quadratic term.

**Starting Point** A neutral way to choose the starting point is to place equal weights on all the c.d.f.s in our basis:  $w_k \equiv 1/K$  ( $k = 1, \dots, K$ ). Then it is necessary to solve the Marčenko-Pastur equation numerically once *before* launching the SLP optimizer, in order to compute the values of  $\check{m}_F(x_j)$  ( $j = 1, \dots, p$ ) that correspond to this initial choice of  $\tilde{H} = \sum_{k=1}^K M_k/K$ . The initial values for  $a_j$  are taken to be  $\text{Re}[\check{m}_F(x_j)]$ , and  $\text{Im}[\check{m}_F(x_j)]$  for  $b_j$  ( $j = 1, \dots, p$ ). If the choice of equal weights  $w_k \equiv 1/K$  for the starting point does not lead to convergence of the optimization program within a pre-specified limit on the maximum number of iterations, we choose random weights  $w_k$  generated i.i.d.  $\sim \text{Uniform}[0,1]$  (rescaled to sum up to one), repeating this process until convergence finally occurs. In the vast majority of cases, the optimization program already converges on the first try. For example, over 1,000 Monte Carlo simulations using the design of Subsection 6.1 with  $p = 100$  and  $n = 300$ , the optimization program converged on the first try 994 times and on the second try the remaining 6 times.

**Optimization Time** Figure 1 gives some information on how the optimization time increases with the matrix dimension.

The main reason for the rate at which the optimization time increases with  $p$  is that the number of grid points in (5.6) increases linearly in  $p$ . This linear rate is not a requirement for our asymptotic results. Therefore, if necessary, it is possible to pick a less-than-linear rate of increase in the number of grid points to speed up the optimization for very large matrices.

**Estimating the Covariance Matrix** Once the SLP optimizer has converged, it generates optimal values  $(a_1^*, \dots, a_p^*)$ ,  $(b_1^*, \dots, b_p^*)$ , and  $(w_1^*, \dots, w_K^*)$ . The first two sets of variables at the optimum are used to estimate the oracle shrinkage factors. From the reconstructed  $\check{m}_F^*(x_j) \equiv a_j^* + ib_j^*$ , we deduce by linear interpolation  $\check{m}_F^*(\lambda_j)$ , for  $j = 1, \dots, p$ . Our estimator

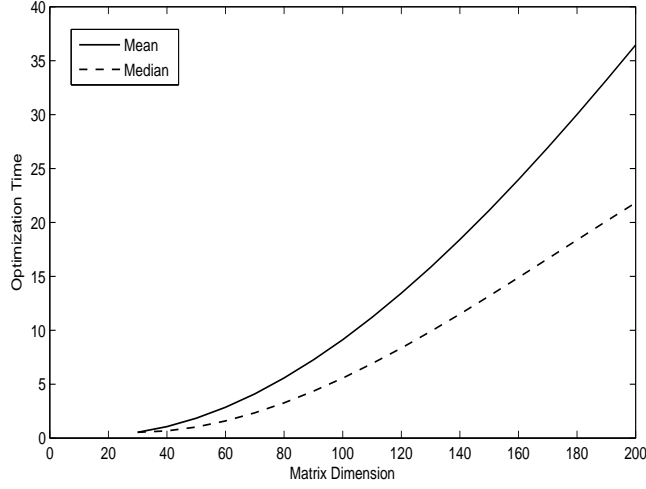


Figure 1: Mean and Median CPU times for Optimization Program as Function of Matrix Dimension. The design is the one of Subsection 6.1 with  $n = 3p$ . Every point is the result of 1,000 Monte Carlo simulations.

of the covariance matrix  $\widehat{S}_n$  is built by keeping the same eigenvectors as the sample covariance matrix, and dividing each sample eigenvalue  $\lambda_j$  by the following correction factor:

$$\left| 1 - \frac{p}{n} - \frac{p}{n} \lambda_j \check{m}_F^*(\lambda_j) \right|^2.$$

Corollary 5.2 assures us that the resulting *bona fide* nonlinear shrinkage estimator is asymptotically equivalent to the oracle estimator  $S_n^{or}$ . Also, we can see that, as the concentration  $\widehat{c}_n = p/n$  gets closer to zero, that is, as we get closer to fixed-dimension asymptotics, the magnitude of the correction becomes smaller. This makes sense because under fixed-dimension asymptotics the sample covariance matrix is a consistent estimator of the population covariance matrix.

**Estimating the Precision Matrix** The output of the same optimization process can also be used to estimate the oracle shrinkage factors for the precision matrix. Our estimator of the precision matrix  $\Sigma_n^{-1}$  is built by keeping the same eigenvectors as the sample covariance matrix, and multiplying the inverse  $\lambda_j^{-1}$  of each sample eigenvalue by the following correction factor:

$$1 - \frac{p}{n} - 2 \frac{p}{n} \lambda_j \operatorname{Re}[\check{m}_F^*(\lambda_j)] .$$

Corollary 5.2 assures us that the resulting *bona fide* nonlinear shrinkage estimator is asymptotically equivalent to the oracle estimator  $P_n^{or}$ .

**Estimating  $H$**  We point out that the optimal values  $(w_1^*, \dots, w_K^*)$  generated from the SLP optimizer yield a consistent estimate of  $H$  in the following fashion:

$$H^* \equiv \sum_{k=1}^K w_k^* M_k .$$

This estimator could be considered an alternative to the estimator introduced by El Karoui (2008). The most salient difference between the two optimization algorithms is that our objective function tries to match  $F_n$  on  $\mathbb{R}$ , whereas his objective function tries to match (a function of)  $m_{F_n}$  on  $\mathbb{C}^+$ . The deeper we go into  $\mathbb{C}^+$ , the more ‘smoothed-out’ is the Stieltjes transform, as it is an analytic function; therefore, the more information is lost. However, the approach of El Karoui (2008) cannot get too close to the real line because  $m_{F_n}$  starts looking like a sum of Dirac functions (which are very ill-behaved) as one gets close to the real line, since  $F_n$  is a step function. In a sense, the approach of El Karoui (2008) is to match a smoothed-out version of a sum of ill-behaved Diracs. In this situation, knowing how much to smooth is rather delicate, and even if it is done well, it still loses information. By contrast, we have no information loss because we operate directly on the real line, and we have no problems with Diracs because we match  $F_n$  instead of its derivative. The price to pay is that our optimization program is not convex, whereas the one of El Karoui (2008) is. But extensive simulations reported in the next section show that off-the-shelf non-convex optimization software — as the commercial package SNOPT — can handle this particular type of a non-convex problem in a fast, robust, and efficient manner.

It would have been of additional interest to compare our estimator of  $H$  to the one of El Karoui (2008) in some simulations. But when we tried to implement his estimator according to the implementation details provided, we were not able to match the results presented in his paper. Furthermore, we were not able to obtain his original software. As a result, we cannot make any definite statements concerning the performance of our estimator of  $H$  compared to the one of El Karoui (2008).

**Remark 5.2** (Cross-Validation Estimator). The implementation of our nonlinear shrinkage estimators is not trivial and also requires the use of a third-party SLP optimizer. It is therefore of interest whether an alternative version exists that is easier to implement and exhibits (nearly) as good finite-sample properties.

To this end an anonymous referee suggested to estimate the quantities  $d_i^*$  of (3.2) by a leave-one-out cross-validation method. In particular, let  $(\lambda_i[k], \dots, \lambda_p[k]); (u_1[k], \dots, u_p[k])$  denote a system of eigenvalues and eigenvectors of the sample covariance matrix computed from all the observed data except for the  $k^{th}$  observation. Then  $d_i^*$  of (3.2) can be approximated by

$$d_i^{cv} \equiv \frac{1}{n} \sum_{k=1}^n (u_i[k]' y_k)^2 ,$$

where the  $p \times 1$  vector  $y_k$  denotes the  $k^{th}$  row of the matrix  $Y_n \equiv X_n \Sigma_n^{1/2}$ .

The motivation here is that

$$(u_i[k]'y_k)^2 = u_i[k]'y_k y_k' u_i[k] ,$$

where  $y_k$  is independent of  $u_i[k]$  and  $\mathbb{E}(y_k y_k') = \Sigma_n$  (even though  $y_k y_k'$  is of rank one only).

We are grateful for this suggestion, since the cross-validation quantities  $d_i^{cv}$  can be computed without the use of any third-party optimization software and the corresponding computer code is very short.

On the other hand, the cross-validation estimator has three disadvantages. First, when  $p$  is large, it takes much longer to compute the cross-validation estimator. The reason is that the spectral decomposition of a  $p \times p$  covariance matrix has to be computed  $n$  times as opposed to only one time. Second, the cross-validation method only applies to the estimation of the covariance matrix  $\Sigma_n$  itself. It is not clear how to adapt this method to the (direct) estimation of the precision matrix  $\Sigma_n^{-1}$  or any other smooth function of  $\Sigma_n$ . Third, the performance of the cross-validation estimator cannot match the performance of our method; see Appendix B. ■

**Remark 5.3.** Another approach proposed recently is the one of Mestre and Lagunas (2006). They use so-called ‘G-estimation’, that is, asymptotic results that assume the sample size  $n$  and the matrix dimension  $p$  go to infinity together, to derive minimum variance beamformers in the context of the spatial filtering of electronic signals. There are several differences between their paper and the present one. First, Mestre and Lagunas (2006) are interested in an optimal  $p \times 1$  weight vector  $w_{opt}$  given by

$$w_{opt} \equiv \underset{w}{\operatorname{argmin}} w' \Sigma_n w , \quad \text{subject to } w' s_d = 1 ,$$

where  $s_d$  is a  $p \times 1$  vector containing signal information. Consequently, Mestre and Lagunas (2006) are ‘only’ interested in a certain functional of  $\Sigma_n$ , while we are interested in the full covariance matrix  $\Sigma_n$  and also in the full precision matrix  $\Sigma_n^{-1}$ . Second, they use the real Stieltjes transform, which is different from the more conventional complex Stieltjes transform used in random matrix theory and in the present paper. Third, their random variables are complex whereas ours are real. The cumulative impact of these differences is best exemplified by the estimation of the precision matrix: Mestre and Lagunas (2006, p.76) recommend  $(1 - p/n)S_n^{-1}$ , which is just a rescaling of the inverse of the sample covariance matrix, whereas our Subsection 3.2 points to a highly nonlinear transformation of the eigenvalues of the sample covariance matrix. ■

## 6 Monte Carlo Simulations

In this section, we present the results of various sets of Monte Carlo simulations designed to illustrate the finite-sample properties of the nonlinear shrinkage estimator  $\widehat{S}_n$ . As detailed in Section 3, the finite-sample optimal estimator in the class of rotation-equivariant estimators is given by  $S_n^*$  as defined in (3.3). Thus, the improvement of the shrinkage estimator  $\widehat{S}_n$  over

the sample covariance matrix will be measured by how closely this estimator approximates  $S_n^*$  relative to the sample covariance matrix. More specifically, we report the Percentage Relative Improvement in Average Loss (PRIAL), which is defined as

$$\text{PRIAL} \equiv \text{PRIAL}(\hat{\Sigma}_n) \equiv 100 \times \left\{ 1 - \frac{\mathbb{E}[\|\hat{\Sigma}_n - S_n^*\|^2]}{\mathbb{E}[\|S_n - S_n^*\|^2]} \right\} \% , \quad (6.1)$$

where  $\hat{\Sigma}_n$  is an arbitrary estimator of  $\Sigma_n$ . By definition, the PRIAL of  $S_n$  is 0% while the PRIAL of  $S_n^*$  is 100%.

Most of the simulations will be designed around a population covariance matrix  $\Sigma_n$  that has 20% of its eigenvalues equal to 1, 40% equal to 3, and 40% equal to 10. This is a particularly interesting and difficult example introduced and analyzed in detail by Bai and Silverstein (1998). For concentration values such as  $c = 1/3$  and below, it displays ‘spectral separation’, that is, the support of the distribution of sample eigenvalues is the union of three disjoint intervals, each one corresponding to a Dirac of population eigenvalues. Detecting this pattern and handling it correctly is a real challenge for any covariance matrix estimation method.

## 6.1 Convergence

The first set of Monte Carlo simulations shows how the nonlinear shrinkage estimator  $\hat{S}_n$  behaves as the matrix dimension  $p$  and the sample size  $n$  go to infinity together. We assume that the concentration ratio  $\hat{c}_n = p/n$  remains constant and equal to  $1/3$ . For every value of  $p$  (and hence  $n$ ), we run 1,000 simulations with normally distributed variables. The PRIAL is plotted in Figure 2. For the sake of comparison, we also report the PRIALs of the oracle  $S_n^{or}$  and the optimal linear shrinkage estimator  $\bar{S}_n$  developed by Ledoit and Wolf (2004).

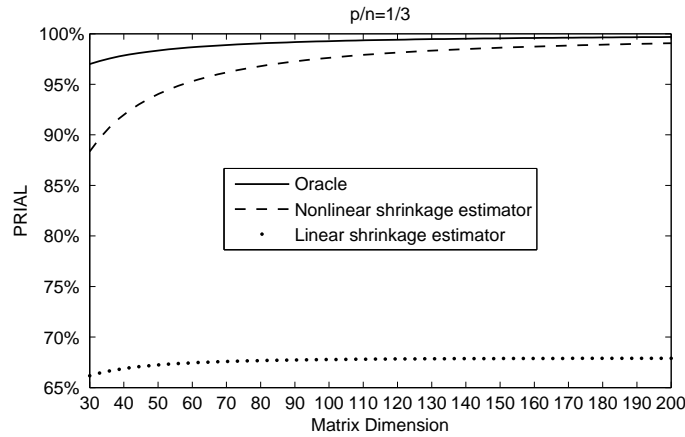


Figure 2: Comparison of the NonLinear vs. Linear Shrinkage Estimators. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. Every point is the result of 1,000 Monte Carlo simulations.

One can see that the performance of the nonlinear shrinkage estimator  $\hat{S}_n$  converges quickly towards that of the oracle and of  $S_n^*$ . Even for relatively small matrices of dimension  $p = 30$ ,

it realizes 88% of the possible gains over the sample covariance matrix. The optimal linear shrinkage estimator  $\bar{S}_n$  performs also well relative to the sample covariance matrix, but the improvement is limited: in general, it does not converge to 100% under large-dimensional asymptotics. This is because there are strong nonlinear effects in the optimal shrinkage of sample eigenvalues. These effects are clearly visible in Figure 3, which plots a typical simulation result for  $p = 100$ .

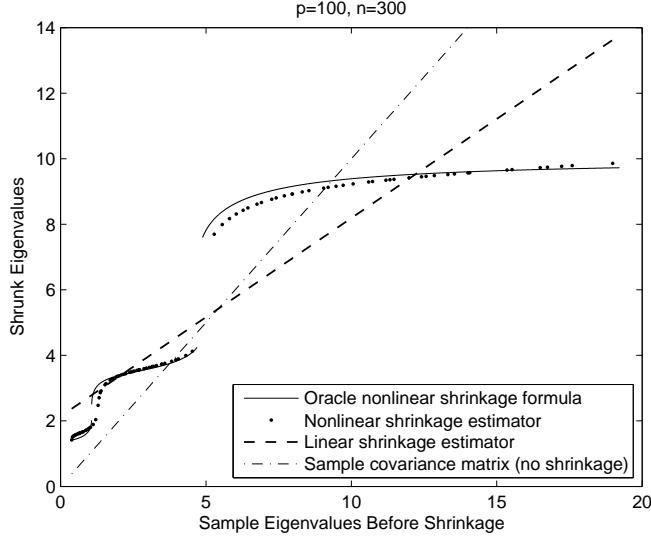


Figure 3: Nonlinearity of the Oracle Shrinkage Formula. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10.  $p = 100$  and  $n = 300$ .

One can see that the nonlinear shrinkage estimator  $\hat{S}_n$  shrinks the eigenvalues of the sample covariance matrix almost as if it ‘knew’ the correct shape of the distribution of population eigenvalues. In particular, the various curves and gaps of the oracle nonlinear shrinkage formula are well picked up and followed by this estimator. By contrast, the linear shrinkage estimator can only use the best linear approximation to this highly nonlinear transformation. We also plot the 45-degrees line as a visual reference to show what would happen if no shrinkage was applied to the sample eigenvalues, that is, if we simply used  $S_n$ .

## 6.2 Concentration

The next set of Monte Carlo simulations shows how the PRIAL of the shrinkage estimators varies as a function of the concentration ratio  $\hat{c}_n = p/n$  if we keep the product  $p \times n$  constant and equal to 9,000. We keep the same population covariance matrix  $\Sigma_n$  as in Subsection 6.1. For every value of  $p/n$ , we run 1,000 simulations with normally distributed variables. The respective PRIALs of  $S_n^{or}$ ,  $\hat{S}_n$ , and  $\bar{S}_n$  are plotted in Figure 4.

One can see that the nonlinear shrinkage estimator performs well across the board, closely in line with the oracle, and always achieves at least 90% of the possible improvement over the sample covariance matrix. By contrast, the linear shrinkage estimator achieves relatively



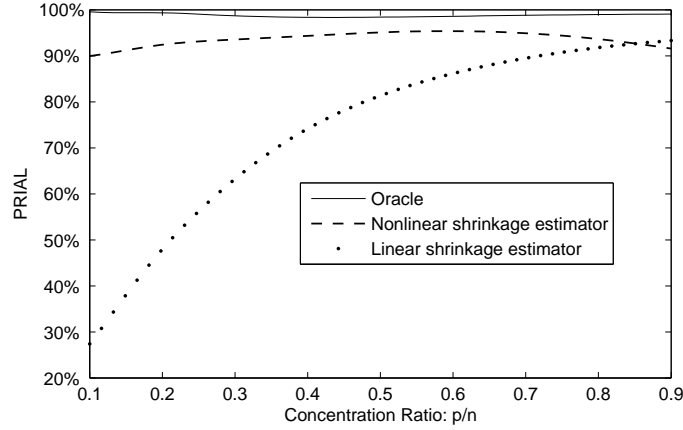


Figure 4: Effect of Varying the Concentration Ratio  $\hat{c}_n = p/n$ . 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. Every point is the result of 1,000 Monte Carlo simulations.

little improvement over the sample covariance matrix when the concentration is low. This is because, when the sample size is large relative to the matrix dimension, there is a lot of precise information about the optimal nonlinear way to shrink the sample eigenvalues that is waiting to be extracted by a suitable nonlinear procedure. By contrast, when the sample size is not so large, the information about the population covariance matrix is relatively fuzzy, therefore a simple linear approximation can achieve up to 93% of the potential gains.

### 6.3 Dispersion

The third set of Monte Carlo simulations shows how the PRIAL of the shrinkage estimators varies as a function of the dispersion of population eigenvalues. We take a population covariance matrix  $\Sigma_n$  with 20% of its eigenvalues equal to 1, 40% equal to  $1 + 2d/9$ , and 40% equal to  $1 + d$ , where the dispersion parameter  $d$  varies from 0 to 20. Thus, for  $d = 0$ ,  $\Sigma_n$  is the identity matrix and, for  $d = 9$ ,  $\Sigma_n$  is the same matrix as in Subsection 6.1. The sample size is  $n = 300$  and the matrix dimension is  $p = 100$ . For every value of  $d$ , we run 1,000 simulations with normally distributed variables. The respective PRIALs of  $S_n^{or}$ ,  $\hat{S}_n$ , and  $\bar{S}_n$  are plotted in Figure 5.

One can see that the linear shrinkage estimator  $\bar{S}_n$  beats the nonlinear shrinkage estimator  $\hat{S}_n$  for very low dispersion levels. For example, when  $d = 0$ , that is, when the population covariance matrix is equal to the identity matrix,  $\bar{S}_n$  realizes 99.9% of the possible improvement over the sample covariance matrix, while  $\hat{S}_n$  realizes ‘only’ 99.4% of the possible improvement. This is because, in this case, linear shrinkage is optimal or (when  $d$  is strictly positive but still small) nearly optimal, hence there is nothing to little to be gained by resorting to a nonlinear shrinkage method. However, as dispersion increases, linear shrinkage delivers less and less improvement over the sample covariance matrix, while nonlinear shrinkage retains a PRIAL above 96%, and close to that of the oracle.

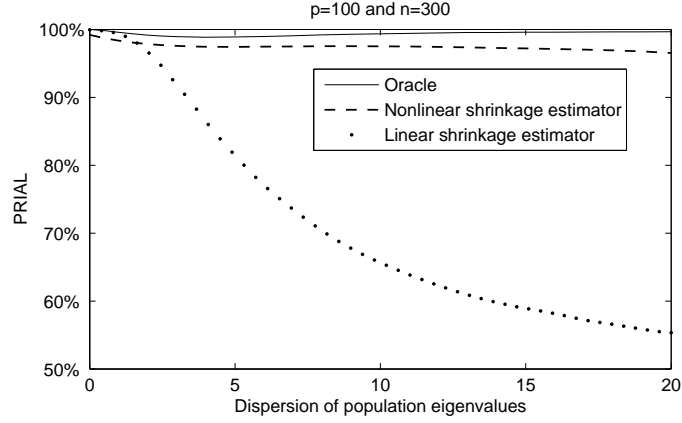


Figure 5: Effect of Varying the Dispersion of Population Eigenvalues. 20% of population eigenvalues are equal to 1, 40% equal to  $1 + 2d/9$ , and 40% equal to  $1 + d$ , where the dispersion parameter  $d$  varies from 0 to 20.  $p = 100$  and  $n = 300$ . Every point is the result of 1,000 Monte Carlo simulations.

#### 6.4 Fat Tails

We also have some results on the effect of non-normality on the performance of the shrinkage estimators. We take the same population covariance matrix as in Subsection 6.1, that is,  $\Sigma_n$  has 20% of its eigenvalues equal to 1, 40% equal to 3, and 40% equal to 10. The sample size is  $n = 300$ , and the matrix dimension is  $p = 100$ . We compare two types of random variates: a Student  $t$  distribution with  $df = 3$  degrees of freedom, and a Student  $t$  distribution with  $df = \infty$  degrees of freedom (which is the Gaussian distribution). For each number of degrees of freedom  $df$ , we run 1,000 simulations. The respective PRIALs of  $S_n^{or}$ ,  $\hat{S}_n$ , and  $\bar{S}_n$  are summarized in Table 1.

	Average Squared Frobenius Loss		PRIAL	
	df = 3	df = $\infty$	df = 3	df = $\infty$
Sample Covariance Matrix	5.856	5.837	0%	0%
Linear Shrinkage Estimator	1.883	1.883	67.84%	67.74%
Nonlinear Shrinkage Estimator	0.128	0.133	97.81%	97.71%
Oracle	0.043	0.041	99.27%	99.30%

Table 1: Effect of Non-normality. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. 1,000 Monte Carlo simulations with  $p = 100$  and  $n = 300$ .

One can see that departure from normality does not have any noticeable effect on performance.

## 6.5 Precision Matrix

The next set of Monte Carlo simulations focuses on estimating the precision matrix  $\Sigma_n^{-1}$ . The definition of the PRIAL, in this subsection only, is given by

$$\text{PRIAL} \equiv \text{PRIAL}(\hat{\Pi}_n) \equiv 100 \times \left\{ 1 - \frac{\mathbb{E} \left[ \|\hat{\Pi}_n - P_n^*\|^2 \right]}{\mathbb{E} \left[ \|S_n^{-1} - P_n^*\|^2 \right]} \right\} \% , \quad (6.2)$$

where  $\hat{\Pi}_n$  is an arbitrary estimator of  $\Sigma_n^{-1}$ . By definition, the PRIAL of  $S_n^{-1}$  is 0% while the PRIAL of  $P_n^*$  is 100%.

We take the same population eigenvalues as in Subsection 6.1. The concentration ratio  $\hat{c}_n = p/n$  is set to the value  $1/3$ . For various values of  $p$  between 30 and 200, we run 1,000 simulations with normally distributed variables. The respective PRIALs of  $P_n^{or}$ ,  $\hat{P}_n$ ,  $\hat{S}_n^{-1}$ , and  $\bar{S}_n^{-1}$  are plotted in Figure 6.

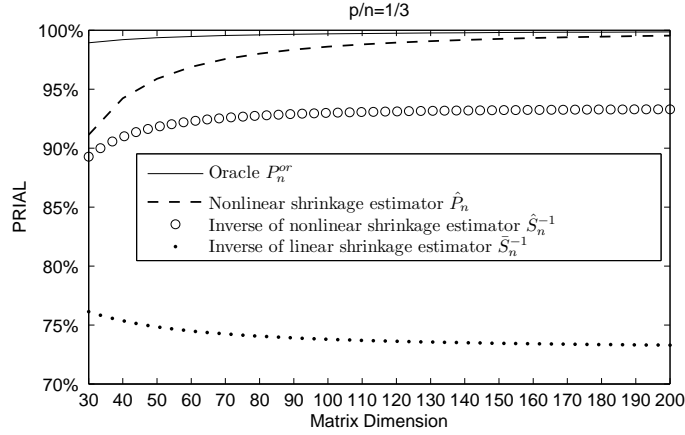


Figure 6: Estimating the Precision Matrix. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. Every point is the result of 1,000 Monte Carlo simulations.

One can see that the nonlinear shrinkage method seems to be just as effective for the purpose of estimating the precision matrix as it is for the purpose of estimating the covariance matrix itself. Moreover, there is a clear benefit in directly estimating the precision matrix by means of  $\hat{P}_n$  as opposed to the indirect estimation by means of  $\hat{S}_n^{-1}$  (which on its own significantly outperforms  $\bar{S}_n^{-1}$ ).

## 6.6 Shape

Next, we study how the nonlinear shrinkage estimator  $\hat{S}_n$  performs for a wide variety of shapes of population spectral densities. This requires using a family of distributions with bounded support and which, for various parameter values, can take on different shapes. The best-suited family for this purpose is the beta distribution. The c.d.f. of the beta distribution with

parameters  $(\alpha, \beta)$  is:

$$\forall x \in [0, 1] \quad F_{(\alpha, \beta)}(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int_0^x t^{\alpha-1}(1-t)^{\beta-1} dt.$$

While the support of the beta distribution is  $[0, 1]$ , we shift it to the interval  $[1, 10]$  by applying a linear transformation. Thanks to the flexibility of the beta family of densities, selecting different parameters  $(\alpha, \beta)$  enables us to generate eight different shapes for the population spectral density: rectangular  $(1, 1)$ , linearly decreasing triangle  $(1, 2)$ , linearly increasing triangle  $(2, 1)$ , circular  $(1.5, 1.5)$ , U-shaped  $(0.5, 0.5)$ , bell-shaped  $(5, 5)$ , left-skewed  $(5, 2)$  and right-skewed  $(2, 5)$ ; see Figure 7 for a graphical illustration.

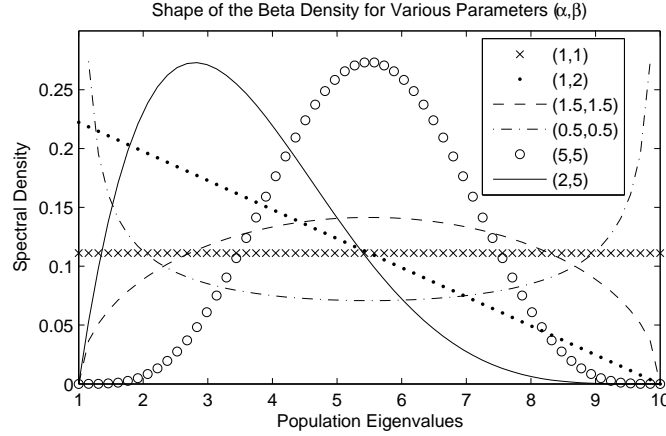


Figure 7: Shape of the Beta Density for Various Parameter Values. The support of the beta density has been shifted to the interval  $[1, 10]$  by a linear transformation. To enhance clarity, the densities corresponding to the parameters  $(2, 1)$  and  $(5, 2)$  have been omitted, since they are symmetric to  $(1, 2)$  and  $(2, 5)$  respectively about the mid-point of the support.

For every one of these eight beta densities, we take the population eigenvalues to be equal to

$$1 + 9 F_{(\alpha, \beta)}^{-1} \left( \frac{i}{p} - \frac{1}{2p} \right), \quad i = 1, \dots, p.$$

The concentration ratio  $\hat{c}_n = p/n$  is equal to  $1/3$ . For various values of  $p$  between 30 and 200, we run 1,000 simulations with normally distributed variables. The PRIAL of the nonlinear shrinkage estimator  $\hat{S}_n$  is plotted in Figure 8.

As in all the other simulations presented above, the PRIAL of the nonlinear shrinkage estimator always exceeds 88%, and more often than not exceeds 95%. To preserve the clarity of the picture, we do not report the PRIALs of the oracle and of the linear shrinkage estimator; but as usual, the nonlinear shrinkage estimator ranked between them.

## 6.7 Fixed-Dimension Asymptotics

Finally, we report a set of Monte-Carlo simulations that departs from the large-dimensional asymptotics assumption under which the nonlinear shrinkage estimator  $\hat{S}_n$  was derived. The

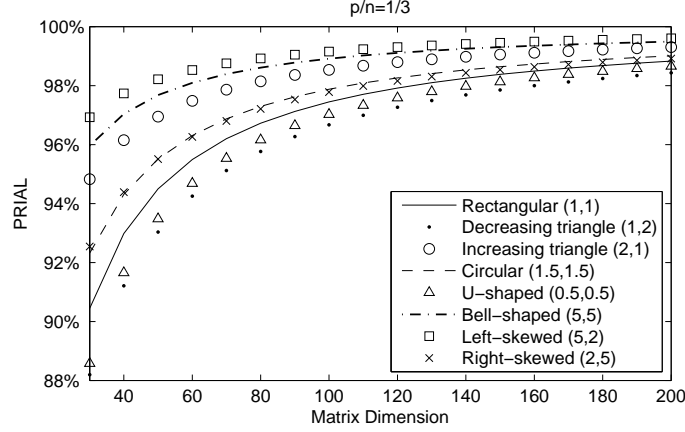


Figure 8: Performance of the Nonlinear Shrinkage with Beta Densities. The various curves correspond to different shapes of the population spectral density. The support of the population spectral density is  $[1, 10]$ .

goal is to compare it against the sample covariance matrix  $S_n$  in the setting where  $S_n$  is known to have certain optimality properties (at least in the normal case): traditional asymptotics, that is, when the number of variables  $p$  remains fixed while the sample size  $n$  goes to infinity. This gives as much advantage to the sample covariance matrix as it can possibly have. We fix the dimension  $p = 100$  and let the sample size  $n$  vary from  $n = 125$  to  $n = 10,000$ . In practice, very few applied researchers are fortunate enough to have as many as  $n = 10,000$  i.i.d. observations, or a concentration ratio  $c = p/n$  as low as 0.01. The respective PRIALs of  $S_n^{or}$ ,  $\hat{S}_n$ , and  $\bar{S}_n$  are plotted in Figure 9.

One crucial difference with all the previous simulations is that the target for the PRIAL is no longer  $S_n^*$ , but instead the population covariance matrix  $\Sigma$  itself, because now  $\Sigma$  can be consistently estimated. Note that, since the matrix dimension is fixed,  $\Sigma_n$  does not change with  $n$ ; therefore, we can drop the subscript  $n$ . Thus, in this subsection only, the definition of the PRIAL is given by

$$\text{PRIAL} \equiv \text{PRIAL}(\hat{\Sigma}_n) \equiv 100 \times \left\{ 1 - \frac{\mathbb{E}[\|\hat{\Sigma}_n - \Sigma\|^2]}{\mathbb{E}[\|S_n - \Sigma\|^2]} \right\} \%,$$

where  $\hat{\Sigma}_n$  is an arbitrary estimator of  $\Sigma$ . By definition, the PRIAL of  $S_n$  is 0% while the PRIAL of  $\Sigma$  is 100%.

In this setting, Ledoit and Wolf (2004) acknowledge that the improvement of the linear shrinkage estimator over the sample covariance matrix vanishes asymptotically, because the optimal linear shrinkage intensity vanishes. Therefore it should be no surprise that the PRIAL of  $\bar{S}_n$  goes to zero in Figure 9. Perhaps more surprising is the continued ability of the oracle and the nonlinear shrinkage estimator to improve by approximately 60% over the sample covariance matrix, even for a sample size as large as  $n = 10,000$ , and with no sign of abating as  $n$  goes to infinity. This is an encouraging result, as our simulation gave every possible advantage to

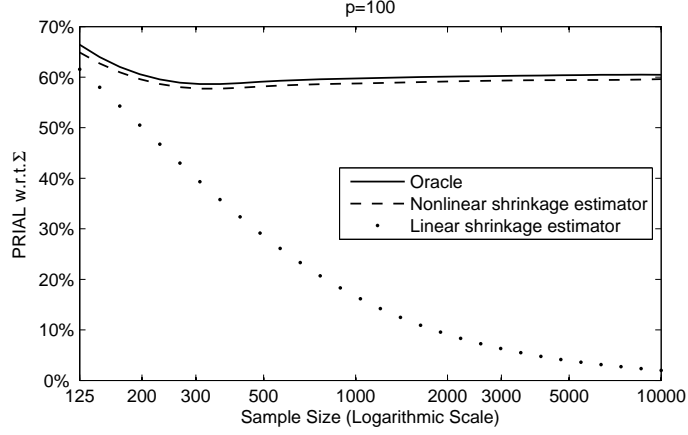


Figure 9: Fixed-Dimension Asymptotics. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. Variables are normally distributed. Every point is the result of 1,000 Monte Carlo simulations.

the sample covariance matrix by placing it in the asymptotic conditions where it possesses well-known optimality properties, and where the earlier linear shrinkage estimator of Ledoit and Wolf (2004) is most disadvantaged.

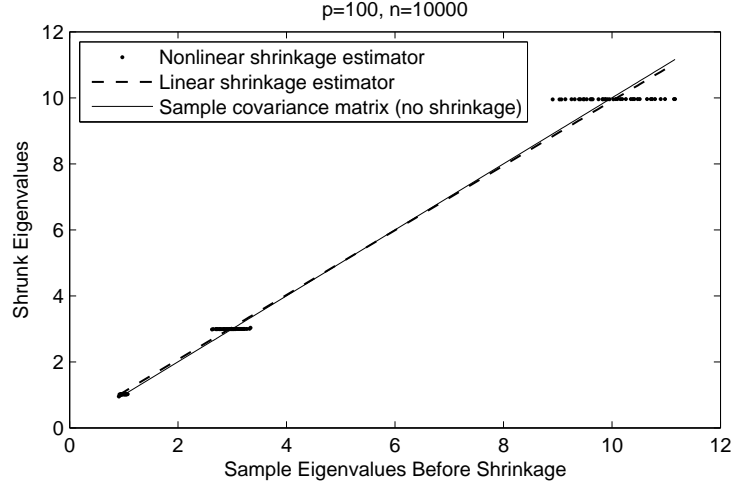


Figure 10: Nonlinear Shrinkage under Fixed-Dimension Asymptotics. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10.  $p = 100$  and  $n = 10,000$ . The oracle is not shown because it is virtually identical to the nonlinear shrinkage estimator.

Intuitively, this is because the oracle shrinkage formula becomes more and more nonlinear as  $n$  goes to infinity for fixed  $p$ . Bai and Silverstein (1998) show that the sample covariance matrix exhibits ‘spectral separation’ when the concentration ratio  $p/n$  is sufficiently small. It means that the sample eigenvalues coalesce into clusters, each cluster corresponding to a Dirac of population eigenvalues. Within a given cluster, the smallest sample eigenvalues need to be nudged upwards, and the largest ones downwards, to the average of the cluster. In other



Monte Carlo studies have confirmed that this estimator yields a sizeable improvement over the indirect method of simply inverting the nonlinear shrinkage estimator of the covariance matrix itself.

The scope of this paper is limited to the case where the matrix dimension is smaller than the sample size. The other case, where the matrix dimension exceeds the sample size, requires certain modifications in the mathematical treatment, and is left for future research.

## References

- Bai, Z. and Silverstein, J. W. (1998). No eigenvalues outside the support of the limiting spectral distribution of large-dimensional random matrices. *Annals of Probability*, 26(1):316–345.
- Bai, Z. D. and Silverstein, J. W. (1999). Exact separation of eigenvalues of large-dimensional sample covariance matrices. *Annals of Probability*, 27(3):1536–1555.
- Bickel, P. J. and Levina, E. (2008). Regularized estimation of large covariance matrices. *Annals of Statistics*, 36(1):199–227.
- Cai, T. and Zhou, H. (2012). Minimax estimation of large covariance matrices under  $\ell_1$  norm. *Statistica Sinica*. Forthcoming.
- El Karoui, N. (2008). Spectrum estimation for large dimensional covariance matrices using random matrix theory. *Annals of Statistics*, 36(6):2757–2790.
- Fan, J., Fan, Y., and Lv, J. (2008). High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics*, 147(1):186–197.
- Geronimo, J. S. and Hill, T. P. (2003). Necessary and sufficient condition that the limit of Stieltjes transforms is a Stieltjes transform. *Journal of Approximation Theory*, 121:54–60.
- Gill, P. E., Murray, W., and Saunders, M. A. (2002). SNOPT: An SQP algorithm for large-scale constrained optimization. *SIAM Journal on Optimization*, 12(4):979–1006.
- Haff, L. R. (1980). Empirical Bayes estimation of the multivariate normal covariance matrix. *Annals of Statistics*, 8:586–597.
- James, W. and Stein, C. (1961). Estimation with quadratic loss. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability* 1, pages 361–380.
- Khare, K. and Rajaratnam, B. (2011). Wishart distributions for decomposable covariance graph models. *Annals of Statistics*, 39(1):514–555.
- Ledoit, O. and Péché, S. (2011). Eigenvectors of some large sample covariance matrix ensembles. *Probability Theory and Related Fields*, 150(1–2):233–264.



- Ledoit, O. and Wolf, M. (2004). A well-conditioned estimator for large-dimensional covariance matrices. *Journal of Multivariate Analysis*, 88(2):365–411.
- Marčenko, V. A. and Pastur, L. A. (1967). Distribution of eigenvalues for some sets of random matrices. *Sbornik: Mathematics*, 1(4):457–483.
- Mestre, X. (2008). On the asymptotic behavior of the sample estimates of eigenvalues and eigenvectors of covariance matrices. *IEEE Transactions on Signal Processing*, 56(11):5353–5368.
- Mestre, X. and Lagunas, M. A. (2006). Finite sample size effect on minimum variance beamformers: Optimum diagonal loading factor for large arrays. *IEEE Transactions on Signal Processing*, 54(1):69–82.
- Perlman, M. D. (2007). *STAT 542: Multivariate Statistical Analysis*. University of Washington (On-Line Class Notes), Seattle, Washington.
- Pollard, D. (1984). *Convergence of Stochastic Processes*. Springer-Verlag, New York.
- Rajaratnam, B., Massam, H., and Carvalho, C. M. (2008). Flexible covariance estimation in graphical Gaussian models. *Annals of Statistics*, 36(6):2818–2849.
- Ravikumar, P., Wawinwright, M., Raskutti, G., and Yu, B. (2008). High-dimensional covariance estimation by minimizing  $\ell_1$ -penalized log-determinant divergence. Technical Report 797, Statistics Department, UC Berkeley.
- Rohde, A. and Tsybakov, A. B. (2010). Estimation of high-dimensional low-rank matrices. *Annals of Statistics*, 39(2):887–930.
- Silverman, B. W. (1986). *Density Estimation for Statistics and Data Analysis*. Chapman & Hall/CRC, Boca Raton.
- Silverstein, J. W. (1995). Strong convergence of the empirical distribution of eigenvalues of large-dimensional random matrices. *Journal of Multivariate Analysis*, 55:331–339.
- Silverstein, J. W. and Choi, S. I. (1995). Analysis of the limiting spectral distribution of large-dimensional random matrices. *Journal of Multivariate Analysis*, 54:295–309.
- Stein, C. (1956). Inadmissibility of the usual estimator for the mean of a multivariate normal distribution. In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, pages 197–206. University of California Press.
- Stein, C. (1975). Estimation of a covariance matrix. Rietz lecture, 39th Annual Meeting IMS. Atlanta, Georgia.
- Won, J.-H., Lim, J., Kim, S.-J., and Rajaratnam, B. (2009). Maximum likelihood covariance estimation with a condition number constraint. Technical Report 2009-10, Department of Statistics, Stanford University.

## A Mathematical Proofs

Before proving Proposition 4.1, it is instructive to first state and prove a simpler result only claiming pointwise convergence of the estimated solutions. We will then see that this simpler proof can be extended relatively easily to also cover the more general claim of uniform convergence.

**Proposition A.1.** *Let  $\{\hat{H}_n\}$  be a sequence of probability measures with  $\hat{H}_n \Rightarrow H$ . Let  $\{\hat{c}_n\}$  be a sequence of positive real numbers with  $\hat{c}_n \rightarrow c$ . Let  $K \subseteq (0, \infty)$  be a compact interval satisfying  $y_x \in K$ . Let  $\hat{y}_{n,x} \equiv \min_{y \in K} g_{\hat{H}_n, \hat{c}_n}(y, x)$ . It then holds that  $\hat{y}_{n,x} \rightarrow y_x$ .*

PROOF. Assume  $K = [k_1, k_2]$ . Define  $B \equiv \{x + iy : x \in [u_1, u_2], y \in K\}$ , which implies  $B \subseteq \mathbb{C}^+$ .

We first claim that

$$m_{L\hat{H}_n}(z) \rightarrow m_{LH}(z) \quad \text{uniformly in } z \in B. \quad (\text{A.1})$$

Recalling that for any c.d.f.  $G$ , we have  $m_{LG}(z) = 1 + z m_G(z)$  and by the compactness of the set  $B$ , this results will follow from

$$m_{\hat{H}_n}(z) \rightarrow m_H(z) \quad \text{uniformly in } z \in B, \quad (\text{A.2})$$

which we establish now.

For fixed  $z \in B$ , consider the function

$$h_z(\tau) \equiv \frac{\tau}{\tau - z}.$$

Then it is easy to see that there exist two finite constants  $d_1, d_2$ , depending only on  $k_1 > 0$  but not on  $z$ , such that

$$|h_z(\tau_1) - h_z(\tau_2)| \leq d_1 |\tau_1 - \tau_2| \quad \text{and} \quad \sup_{\tau} |h_z(\tau)| \leq d_2. \quad (\text{A.3})$$

The fact that convergence in distribution of  $\hat{H}_n$  to  $H$  is equivalent to convergence to zero of the bounded-Lipschitz metric between  $\hat{H}_n$  and  $H$  then implies (A.2); for example, see Pollard (1984, Example 22). In turn, we have thus established (A.1) as well. But (A.1) immediately implies

$$g_{\hat{H}_n, \hat{c}_n}(y, x) \rightarrow g_{H, c}(y, x) \quad \text{uniformly in } y \in K. \quad (\text{A.4})$$

We note the following two facts:

$$\forall \varepsilon > 0 \quad \exists \delta > 0 \text{ such that } \inf_{y \in K, |y - y_x| \geq \varepsilon} g_{H, c}(y, x) \geq \delta \quad (\text{A.5})$$

and

$$g_{\hat{H}_n, \hat{c}_n}(\hat{y}_{n,x}, x) = o(1), \quad (\text{A.6})$$

where (A.6) follows from  $g_{\hat{H}_n, \hat{c}_n}(\hat{y}_{n,x}, x) \leq g_{\hat{H}_n, \hat{c}_n}(y_x, x)$ , (A.4), and  $g_{H, c}(y_x, x) = 0$ .

By the triangular inequality,

$$\begin{aligned}
g_{H,c}(\widehat{y}_{n,x}, x) &\leq |g_{H,c}(\widehat{y}_{n,x}) - g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x})| + |g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x})| \\
&= |g_{H,c}(\widehat{y}_{n,x}) - g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x})| + o(1) \quad \text{by (A.6)} \\
&= o(1) + o(1) \quad \text{by (A.4)} \\
&= o(1) .
\end{aligned}$$

This last result together with (A.5) now implies  $\widehat{y}_{n,x} \rightarrow y_x$ . ■

PROOF OF PROPOSITION 4.1. We start with part (i). Assume  $K = [k_1, k_2]$ . Define  $B \equiv \{x + i y : x \in [u_1, u_2], y \in K\}$ , which implies  $B \subseteq \mathbb{C}^+$ .

By the same arguments leading up to (A.4) we can more generally establish that

$$g_{\widehat{H}_n, \widehat{c}_n}(z) \rightarrow g_{H,c}(z) \quad \text{uniformly in } z \in B . \quad (\text{A.7})$$

We note the following two facts:

$$\forall \varepsilon > 0 \quad \exists \delta > 0 \text{ such that } \inf_{x \in [u_1 + \eta, u_2 - \eta]} \left\{ \inf_{y \in K, |y - y_x| \geq \varepsilon} g_{H,c}(y, x) \right\} \geq \delta \quad (\text{A.8})$$

and

$$\sup_{x \in [u_1 + \eta, u_2 - \eta]} g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x}, x) = o(1) , \quad (\text{A.9})$$

where (A.9) follows from  $g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x}, x) \leq g_{\widehat{H}_n, \widehat{c}_n}(y_x, x)$ , (A.7), and  $g_{H,c}(y_x, x) = 0$ .

To simplify the notation, let  $I \equiv [u_1 + \eta, u_2 - \eta]$ . By the triangular inequality,

$$\begin{aligned}
\sup_{x \in I} g_{H,c}(\widehat{y}_{n,x}, x) &\leq \sup_{x \in I} |g_{H,c}(\widehat{y}_{n,x}) - g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x})| + \sup_{x \in I} |g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x})| \\
&= \sup_{x \in I} |g_{H,c}(\widehat{y}_{n,x}) - g_{\widehat{H}_n, \widehat{c}_n}(\widehat{y}_{n,x})| + o(1) \quad \text{by (A.9)} \\
&= o(1) + o(1) \quad \text{by (A.7)} \\
&= o(1) .
\end{aligned}$$

This last result together with (A.9) now implies  $\widehat{y}_{n,x} \rightarrow y_x$  uniformly in  $x \in I = [u_1 + \eta, u_2 - \eta]$ .

Part (ii) is proven analogously to part (i) by restricting attention to the set of probability one on which  $\widehat{H}_n \Rightarrow H$  happens. ■

PROOF OF PROPOSITION 4.2. The proof is similar to the proof of Proposition 4.1. The details are left to the reader. ■

PROOF OF PROPOSITION 4.3. We start with part (i)(a). Fix  $\lambda \in [\widetilde{z}_1 + \widetilde{\delta}, \widetilde{z}_2 - \widetilde{\delta}]$ . Consider

$$\left| \check{m}_{F_{\widehat{H}_n, \widehat{c}_n}}(\lambda) - \check{m}_F(\lambda) \right| = \left| \frac{1 - \widehat{c}_n}{\widehat{c}_n \lambda} - \frac{1}{\widehat{c}_n} \frac{1}{\widehat{v}_{n,\lambda}} - \left( \frac{1 - c}{c \lambda_x} - \frac{1}{c} \frac{1}{v_\lambda} \right) \right| .$$

The function mapping  $\lambda$  onto  $v_\lambda$  is continuous, and therefore uniformly continuous, in  $\lambda \in [\tilde{z}_1, \tilde{z}_2]$ . As  $\lambda$  varies in  $[\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ , the resulting  $v_\lambda$  varies in a compact region in  $\mathbb{C}^+$ . Therefore, for any  $\xi > 0$ , there exists  $\kappa > 0$  such that

$$|\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda) - \check{m}_F(\lambda)| < \xi \quad \text{as long as} \quad \max\{|\hat{c}_n - c|, |\hat{v}_{n, \lambda} - v_\lambda|\} < \kappa.$$

First, we can find  $N_1$  such that  $|\hat{c}_n - c| < \kappa$  for all  $n \geq N_1$ . Second, by part (i) of Proposition 4.2, we can find  $N_2$  such that  $|\hat{v}_{n, \lambda} - v_\lambda| < \kappa$  for all  $n \geq N_2$ , uniformly in  $\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]$ . Define  $N \equiv \max\{N_1, N_2\}$ . Then for all  $n \geq N$ , it holds that

$$|\check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda) - \check{m}_F(\lambda)| < \xi, \quad \text{uniformly in } \lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}].$$

Since  $\xi$  can be chosen arbitrarily small, part (i)(a) obtains.

We now turn to part (i)(b). For any  $\tilde{\delta} > 0$ , it holds

$$\begin{aligned} \|\hat{S}_n - S_n^{or}\|^2 &= \frac{1}{p} \sum_{i=1}^p \left( \frac{\lambda_i}{|1 - \hat{c}_n - \hat{c}_n \lambda_i \check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda_i)|^2} - \frac{\lambda_i}{|1 - c - c \lambda_i \check{m}_F(\lambda_i)|^2} \right)^2 \\ &= \frac{1}{p} \sum_{\lambda_i \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]} \left( \frac{\lambda_i}{|1 - \hat{c}_n - \hat{c}_n \lambda_i \check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda_i)|^2} - \frac{\lambda_i}{|1 - c - c \lambda_i \check{m}_F(\lambda_i)|^2} \right)^2 \\ &\quad + \frac{1}{p} \sum_{\lambda_i \notin [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]} \left( \frac{\lambda_i}{|1 - \hat{c}_n - \hat{c}_n \lambda_i \check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda_i)|^2} - \frac{\lambda_i}{|1 - c - c \lambda_i \check{m}_F(\lambda_i)|^2} \right)^2 \\ &\equiv A + B. \end{aligned}$$

By our general set of assumptions, in particular Assumption (A4), combined with the results of Bai and Silverstein (1998) and Mestre (2008, Section II), there exist two finite, non-zero constants  $\kappa_1 < \kappa_2$  such that  $\kappa_1 \leq \lambda_i \leq \kappa_2$  for all  $i = 1, \dots, p$  and for all  $n$  large enough.

Fix  $\varepsilon > 0$ . First, we can pick  $\tilde{\delta}$  small enough to achieve  $B \leq \varepsilon/2$  eventually. To appreciate why, denote by  $\mu(\tilde{\delta})$  the mass that  $F$  assigns to the set  $[\tilde{z}_1, \tilde{z}_1 + \tilde{\delta}] \cup [\tilde{z}_2 - \tilde{\delta}, \tilde{z}_2]$ , satisfying  $\mu(\tilde{\delta}) \rightarrow 0$  as  $\tilde{\delta} \rightarrow 0$ . Then it is not too difficult to see that there exists a finite constant  $\Delta$ , possibly depending on  $H$  and  $c$ , such that  $B \leq \Delta \mu(\tilde{\delta})$ , for  $n$  sufficiently large. The reason, in addition to  $\kappa_1 \leq \lambda_i \leq \kappa_2$ , is that also the correction factors  $|1 - \hat{c}_n - \hat{c}_n \lambda_i \check{m}_{F_{\hat{H}_n, \hat{c}_n}}(\lambda_i)|^2$  and  $1/|1 - c - c \lambda_i \check{m}_F(\lambda_i)|^2$  are bounded away from infinity. Then, choose  $\tilde{\delta}$  small enough so that  $\mu(\tilde{\delta}) \leq (2/\varepsilon)\Delta$ .

Having chosen and fixed  $\tilde{\delta}$ , the first half of the assertion ensures that  $A \leq \varepsilon/2$  eventually. Again, we use here that  $\kappa_1 \leq \lambda_i \leq \kappa_2$  and that also the correction factors  $1/|1 - c - c \lambda_i \check{m}_F(\lambda_i)|^2$  are bounded away from infinity. This demonstrates part (i)(b).

Part(i)(c) can be handled in a very similar fashion.

Part (ii) is proven analogously to part (i) by focusing on the set of probability one on which  $\hat{H}_n \Rightarrow H$  happens. ■

Before proving Theorem 5.1, we need to establish some auxiliary results.

Recall the following notation. For a grid  $Q$  on the real line and for two c.d.f.s  $G_1$  and  $G_2$ , define

$$||G_1 - G_2||_Q \equiv \sup_{t \in Q} |G_1(t) - G_2(t)| .$$

**Lemma A.1.** *Let  $\{G_n\}$  and  $G$  be c.d.f.s on the real line, with the support of  $G$  being compact. Let  $\{Q_n\}$  be a sequence of grids on the real line, asymptotically covering the support of  $G$ , with grid sizes  $\{\gamma_n\}$  satisfying  $\gamma_n \rightarrow 0$ .*

*If  $G$  is continuous, then  $G_n \Rightarrow G$ . In particular,  $\sup_t |G_n(t) - G(t)| \rightarrow 0$ .*

PROOF. Denote the compact support of  $G$  by  $[a, b]$ . To prove the first part of the assertion, let  $\varepsilon > 0$ . Fix  $\delta > 0$  such that for all  $t < t'$  with  $t' - t < \delta$ , it holds  $G(t') - G(t) < \varepsilon/4$ . Also fix  $\phi > 0$ . First, there exists  $N_1$  such that  $\gamma_n < \delta$  for all  $n \geq N_1$ . Second, there exists  $N_2$  such that  $\sup_{t \in Q_n} |G_n(t) - G(t)| < \varepsilon/4$  for all  $n \geq N_2$ . Third, there exists  $N_3$  such that  $Q_n$  covers  $[a + \phi, b - \phi]$  for all  $n \geq N_3$ . Set  $N \equiv \max\{N_1, N_2, N_3\}$ . For an arbitrary  $t \in [a + \phi, b - \phi]$  and for  $n \geq N$ , let  $t_n \equiv \max\{\tilde{t} : \tilde{t} \in Q_n, \tilde{t} \leq t\}$  and  $t'_n \equiv \min\{\tilde{t} : \tilde{t} \in Q_n, \tilde{t} \geq t\}$ , which implies  $t_n - t'_n < \delta$ . Then, for all  $n \geq N$ ,

$$\begin{aligned} |G_n(t) - G(t)| &\leq |G_n(t_n) - G(t'_n)| + |G_n(t'_n) - G(t_n)| \\ &\leq |G_n(t_n) - G(t_n)| + |G_n(t'_n) - G(t'_n)| + \frac{\varepsilon}{4} + \frac{\varepsilon}{4} \\ &\leq \frac{\varepsilon}{4} + \frac{\varepsilon}{4} + \frac{\varepsilon}{4} + \frac{\varepsilon}{4} = \varepsilon . \end{aligned}$$

Therefore,  $G_n(t)$  converges to  $G(t)$  for all  $t \in [a + \phi, b - \phi]$ ; and since  $\phi$  can be chosen arbitrarily,  $G_n(t)$  converges to  $G(t)$  for all  $t \in (a, b)$ . By picking  $\phi$  sufficiently small such that  $|G(a + \phi)| \leq \varepsilon$  and  $|G(b - \varepsilon)| \geq 1 - \varepsilon$ , and by the monotonicity of c.d.f.s, it also follows that  $|G_n(t)| \leq 2\varepsilon$  for all  $t \leq a$  as well as  $|G_n(t)| \geq 1 - 2\varepsilon$  for all  $t \geq b$  as long as  $n \geq N$  (where  $N$  of course is allowed to depend on  $\phi$ .) Therefore,  $G_n(t)$  converges to  $G(t)$  for all  $t$ , which establishes  $G_n \Rightarrow G$ . The second part of the assertion follows immediately from the first part and Polya's Theorem. ■

**Lemma A.2.** *Let  $G$  be a probability measure with compact support contained in  $(0, +\infty)$  and let  $d > 0$ . Let  $\{\hat{G}_n\}$  be a sequence of probability measures on the nonnegative real line with  $\hat{G}_n \Rightarrow G$  and let  $\{\hat{d}_n\}$  be a sequence of positive real numbers with  $\hat{d}_n \rightarrow d$ . Also assume that there exists an interval  $[a, b]$  contained in  $(0, +\infty)$  such that  $\text{Supp}(\hat{G}_n) \subseteq [a, b]$  for all  $n$  large enough.*

*Then  $F_{\hat{G}_n, \hat{d}_n} \Rightarrow F_{G, d}$ .*

PROOF. Let  $z_j \equiv i \cdot (1 + 1/j)$ , for  $j = 1, 2, \dots$ . Then  $\{z_j\}$  is an infinite sequence in  $\mathbb{C}^+$  with limit point  $z_0 \equiv i \in \mathbb{C}^+$ . By Theorem 2 of Geronimo and Hill (2003), it is sufficient to show that, for all  $z_j$ ,

$$m_{F_{\hat{G}_n, \hat{d}_n}}(z_j) \rightarrow m_{F_{G, d}}(z_j) . \quad (\text{A.10})$$

Recall the notation  $m_{F_{\tilde{H}, \tilde{c}}}$  for the solution of the Marčenko-Pastur equation, for any probability measure  $\tilde{H}$  and for any  $\tilde{c} > 0$ . Namely, for each  $z \in \mathbb{C}^+$ ,  $m_{F_{\tilde{H}, \tilde{c}}}(z)$  is the unique solution

for  $m \in \mathbb{C}^+$  to the equation

$$m = \int_{-\infty}^{+\infty} \frac{1}{\tau [1 - \tilde{c} - \tilde{c} z m] - z} d\tilde{H}(\tau) .$$

Also, define the function

$$\forall m, z \in \mathbb{C} \quad f_{\tilde{H}, \tilde{c}}(m, z) \equiv \left| m - \int_{-\infty}^{+\infty} \frac{1}{\tau [1 - \tilde{c} - \tilde{c} z m] - z} d\tilde{H}(\tau) \right| .$$

In this notation, for a given  $z \in \mathbb{C}^+$ ,  $m_{F_{\tilde{H}, \tilde{c}}}(z)$  is the unique solution for  $m \in \mathbb{C}^+$  to the equation  $f_{\tilde{H}, \tilde{c}}(m, z) = 0$ . Alternatively,  $m_{F_{\tilde{H}, \tilde{c}}}(z)$  is the unique minimizer over  $m \in \mathbb{C}^+$  of the function  $f_{\tilde{H}, \tilde{c}}(\cdot, z)$ . Note that the Stieltjes transform of any probability measure maps  $\mathbb{C}^+$  onto  $\mathbb{C}^+$ . So if  $z \in \mathbb{C}^+$ , then  $m_{F_{\tilde{H}, \tilde{c}}}(z)$  is actually the unique minimizer over  $m \in \mathbb{C}$  of the function  $f_{\tilde{H}, \tilde{c}}(\cdot, z)$ .

Fix  $z_j$  and use the following abbreviations:  $\hat{m}_{n, z_j} \equiv m_{F_{\hat{G}_n, \hat{d}_n}}(z_j)$  and  $m_{z_j} \equiv m_{F_{G, d}}(z_j)$ . The goal then is to show that  $\hat{m}_{n, z_j} \rightarrow m_{z_j}$ .

We claim that there exists a compact set  $S \subseteq \mathbb{C}$  such that  $\hat{m}_{n, z_j} \in S$  for all  $n$ . The proof is by means of contradiction. Assume the claim does not hold. Then there exists a subsequence  $\{n_k\}$  such that  $|\hat{m}_{n_k, z_j}| \rightarrow \infty$ . By the combined assumptions, we can then find  $\Delta > 0$  such that, for all  $n_k$  large enough and for all  $\tau \in [a, b]$ ,

$$\frac{1}{|\tau [1 - \hat{d}_{n_k} - \hat{d}_{n_k} z_j \hat{m}_{n_k, z_j}] - z_j|} \leq \Delta$$

implying that for all  $n_k$  large enough,

$$\begin{aligned} |\hat{m}_{n_k, z_j}| &= \left| \int_{-\infty}^{+\infty} \frac{1}{\tau [1 - \hat{d}_{n_k} - \hat{d}_{n_k} z_j \hat{m}_{n_k, z_j}] - z_j} d\hat{G}_{n_k}(\tau) \right| \\ &= \left| \int_a^b \frac{1}{\tau [1 - \hat{d}_{n_k} - \hat{d}_{n_k} z_j \hat{m}_{n_k, z_j}] - z_j} d\hat{G}_{n_k}(\tau) \right| \\ &\leq \int_a^b \frac{1}{|\tau [1 - \hat{d}_{n_k} - \hat{d}_{n_k} z_j \hat{m}_{n_k, z_j}] - z_j|} d\hat{G}_{n_k}(\tau) \\ &\leq (b - a) \Delta . \end{aligned}$$

But this is in contradiction to  $|\hat{m}_{n_k, z_j}| \rightarrow \infty$ . We may assume w.l.o.g. that  $m_{z_j} \in S$  as well; otherwise sufficiently enlarge  $S$ .

We may further assume that  $S$  is ‘doubly nonnegative’, that is, for all  $m \in S$ , it holds that  $\operatorname{Re}(m) \geq 0$  as well as  $\operatorname{Im}(m) \geq 0$ . The reason is as follows. On the one hand,  $\operatorname{Re}(\hat{m}_{n, z_j}) \geq 0$  for all  $n$  as well as  $\operatorname{Re}(m_{z_j}) \geq 0$ . For example, recalling that  $\operatorname{Re}(z_j) = 0$ ,

$$\operatorname{Re}(m_{z_j}) = \operatorname{Re}(m_{F_{G, d}}(z_j)) = \int_{-\infty}^{\infty} \operatorname{Re} \left( \frac{1}{\lambda - z_j} \right) dF_{G, d}(\lambda) = \int_{-\infty}^{\infty} \frac{\lambda}{|\lambda - z_j|^2} dF_{G, d}(\lambda) ,$$

where  $F_{G, d}$  places all its mass on  $[0, +\infty)$ . On the other hand, since  $z_j \in \mathbb{C}^+$ , also  $\operatorname{Im}(\hat{m}_{n, z_j}) > 0$  for all  $n$  as well as  $\operatorname{Im}(m_{z_j}) > 0$ .

We next claim that

$$f_{\widehat{G}_n, d}(m, z_j) \rightarrow f_{G, d}(m, z_j) \quad \text{uniformly in } m \in S. \quad (\text{A.11})$$

To see why, for  $m \in S$ , consider the function

$$h_{m, z_j}(\tau) \equiv \frac{1}{\tau [1 - d - d z_j m] - z_j}.$$

Since  $S$  is compact,  $\min\{\operatorname{Re}(m), \operatorname{Im}(m)\} \geq 0$  for all  $m \in S$ ,  $\operatorname{Re}(z_j) = 0$ , and  $\operatorname{Im}(z_j) \geq 1$ , there exist two finite constants  $d_1$  and  $d_2$ , allowed to depend on  $S$ , such that

$$|h_{m, z_j}(\tau_1) - h_{m, z_j}(\tau_2)| \leq d_1 |\tau_1 - \tau_2| \quad \text{for } \tau_1, \tau_2 \in [0, +\infty) \quad (\text{A.12})$$

and

$$\sup_{\tau \in [0, +\infty)} |h_{m, z_j}(\tau)| \leq d_2. \quad (\text{A.13})$$

To see why, start with (A.13). It holds that

$$\operatorname{Im}(\tau [1 - d - d z_j m] - z_j) = -(\tau d [\operatorname{Re}(z_j) \operatorname{Im}(m) + \operatorname{Im}(z_j) \operatorname{Re}(m)] + \operatorname{Im}(z_j)).$$

Under the stated conditions,  $\operatorname{Re}(z_j) \operatorname{Im}(m) + \operatorname{Im}(z_j) \operatorname{Re}(m) \geq 0$  and  $\operatorname{Im}(z_j) \geq 1$ . Therefore, as long as  $\tau \geq 0$ , it follows that

$$|\tau [1 - d - d z_j m] - z_j| \geq |\operatorname{Im}(\tau [1 - d - d z_j m] - z_j)| \geq 1,$$

implying that we may choose  $d_2 \equiv 1$ .

Moving on to (A.12), let  $\Delta \equiv \max_{m \in S} |m|$  and note that  $|z_j| \leq 2$ . Therefore, for any  $\tau_1, \tau_2 \in [0, +\infty)$ ,

$$\begin{aligned} |h_{m, z_j}(\tau_1) - h_{m, z_j}(\tau_2)| &= |\tau_1 - \tau_2| \left| \frac{1 - d - d z_j m}{(\tau_1 [1 - d - d z_j m] - z_j)(\tau_2 [1 - d - d z_j m] - z_j)} \right| \\ &= |\tau_1 - \tau_2| \frac{|1 - d - d z_j m|}{|\tau_1 [1 - d - d z_j m] - z_j| |\tau_2 [1 - d - d z_j m] - z_j|} \\ &= |\tau_1 - \tau_2| \frac{|1 - d - d z_j m|}{|\tau_1 [1 - d - d z_j m] - z_j| |\tau_2 [1 - d - d z_j m] - z_j|} \\ &\leq |\tau_1 - \tau_2| (1 + d + 2 d \Delta), \end{aligned}$$

implying that we may choose  $d_1 \equiv (1 + d + 2 d \Delta)$ .

Recall that convergence in distribution of  $\widehat{G}_n$  to  $G$  is equivalent to convergence to zero of the bounded-Lipschitz metric between  $\widehat{G}_n$  and  $G$ ; for example, see Pollard (1984, Example 22). Furthermore, since  $\widehat{G}_n$  and  $G$  put all their mass on  $[0, \infty)$ , it is actually sufficient to start all

integrals at  $\tau = 0$  rather than at  $\tau = -\infty$ . Therefore,

$$\begin{aligned}
\int_{-\infty}^{+\infty} \frac{d\widehat{G}_n(\tau)}{\tau [1 - d - d z_j m] - z_j} &= \int_0^{+\infty} \frac{1}{\tau [1 - d - d z_j m] - z_j} d\widehat{G}_n(\tau) \\
&= \int_0^{\infty} h_{m, z_j}(\tau) d\widehat{G}_n(\tau) \\
&\rightarrow \int_0^{\infty} h_{m, z_j}(\tau) dG(\tau) \\
&= \int_0^{+\infty} \frac{1}{\tau [1 - d - d z_j m] - z_j} dG(\tau) \\
&= \int_{-\infty}^{+\infty} \frac{1}{\tau [1 - d - d z_j m] - z_j} dG(\tau) \quad \text{uniformly in } m \in S,
\end{aligned}$$

which establishes (A.11). But (A.11), combined with the compactness of  $S$ , further implies that also

$$f_{\widehat{G}_n, \widehat{d}_n}(m, z_j) \rightarrow f_{G, d}(m, z_j) \quad \text{uniformly in } m \in S. \quad (\text{A.14})$$

Summing up, we have the following facts: First, there exists a compact set  $S \subseteq C$  such that  $\widehat{m}_{n, z_j}$  is the unique minimizer of  $f_{\widehat{G}_n, \widehat{d}_n}(\cdot, z_j)$  over  $m \in S$  and  $m_{z_j}$  is the unique minimizer of  $f_{G, d}(\cdot, z_j)$  over  $m \in S$ . Second, the function  $f_{G, d}(\cdot, z_j)$  is continuous in  $m$ . Third, the uniform convergence (A.14).

With these facts,  $\widehat{m}_{n, z_j} \rightarrow m_{z_j}$  follows from arguments very similar to those used in the proof of Proposition A.1. ■

**PROOF OF THEOREM 5.1.** We start with the proof of part (i). Since  $c < 1$ , it follows from Silverstein and Choi (1995) that  $F$  is continuously differentiable on all of  $\mathbb{R}$ . By Polya's Theorem it then follows that  $\sup_t |F_n(t) - F(t)| \rightarrow 0$  a.s., implying that  $\|F_n - F\|_{Q_n} \rightarrow 0$  a.s. Also, by construction,  $\|F_{\widehat{H}_n, \widehat{c}_n} - F_n\|_{Q_n} \leq \|F_{H, \widehat{c}_n} - F_n\|_{Q_n}$ . Therefore,

$$\begin{aligned}
\|F_{\widehat{H}_n, \widehat{c}_n} - F\|_{Q_n} &\leq \|F_{\widehat{H}_n, \widehat{c}_n} - F_n\|_{Q_n} + \|F_n - F\|_{Q_n} \\
&\leq \|F_{H, \widehat{c}_n} - F_n\|_{Q_n} + \|F_n - F\|_{Q_n} \\
&\leq \|F_{H, \widehat{c}_n} - F_{H, c}\|_{Q_n} + \|F_{H, c} - F_n\|_{Q_n} + \|F_n - F\|_{Q_n} \\
&= \|F_{H, \widehat{c}_n} - F\|_{Q_n} + 2\|F_n - F\|_{Q_n} \rightarrow 0 \quad \text{a.s.},
\end{aligned}$$

where Lemma A.2 in conjunction with Polya's Theorem is used to show that  $\|F_{H, \widehat{c}_n} - F\|_{Q_n} \rightarrow 0$ . The desired result now follows from Lemma A.1.

We now turn to proving part (ii). By Theorem 2 of Geronimo and Hill (2003), it is sufficient to show that there exists an infinite sequence  $\{v_j\}$  in  $\mathbb{C}^+$  with a limit point  $v_0 \in \mathbb{C}^+$  such that

$$m_{\widehat{H}_n}(v_j) \rightarrow m_H(v_j) \quad \text{a.s. } \forall j. \quad (\text{A.15})$$

Recall the notation  $m_{F_{\widehat{H}, \widehat{c}}}$  for the solution of the Marčenko-Pastur equation, for any probability measure  $\widehat{H}$  and for any  $\widehat{c} > 0$ . Namely, for each  $z \in \mathbb{C}^+$ ,  $m_{F_{\widehat{H}, \widehat{c}}}(z)$  is the unique solution



for  $m \in \mathbb{C}^+$  to the equation

$$m = \int_{-\infty}^{+\infty} \frac{1}{\tau [1 - \tilde{c} - \tilde{c} z m] - z} d\tilde{H}(\tau) .$$

Analogously, to Subsection 2.2, also let

$$\forall x \in \mathbb{R} \quad \underline{F}_{\tilde{H}, \tilde{c}}(x) \equiv (1 - \tilde{c}) \mathbb{1}_{[0, +\infty)}(x) + \tilde{c} F_{\tilde{H}, \tilde{c}}(x)$$

and

$$\forall z \in \mathbb{C}^+ \quad m_{\underline{F}_{\tilde{H}, \tilde{c}}}(z) \equiv \frac{\tilde{c} - 1}{z} + \tilde{c} m_{F_{\tilde{H}, \tilde{c}}}(z) .$$

Hence, for each  $z \in \mathbb{C}^+$ ,  $m_{\underline{F}_{\tilde{H}, \tilde{c}}}(z)$  is the unique solution for  $m \in \mathbb{C}^+$  to the equation

$$m = - \left[ z - \tilde{c} \int_{-\infty}^{+\infty} \frac{\tau}{1 + \tau m} d\tilde{H}(\tau) \right]^{-1} .$$

On  $\mathbb{C}^+$ ,  $m_{\underline{F}_{\tilde{H}, \tilde{c}}}(z)$  has a unique inverse, given by

$$\forall m \in m_{\underline{F}_{\tilde{H}, \tilde{c}}}(\mathbb{C}^+) \quad z_{\underline{F}_{\tilde{H}, \tilde{c}}}(m) \equiv -\frac{1}{m} + \tilde{c} \int_{-\infty}^{+\infty} \frac{\tau}{1 + \tau m} d\tilde{H}(\tau) .$$

Note that both  $m_{\underline{F}_{\tilde{H}, \tilde{c}}}$  and  $z_{\underline{F}_{\tilde{H}, \tilde{c}}}$  are continuous functions. Also in this notation, we have  $\underline{F} = \underline{F}_{H, c}$ ,  $m_{\underline{F}} = m_{\underline{F}_{H, c}}$ , and  $z_{\underline{F}} = z_{\underline{F}_{H, c}}$  then.

As Silverstein and Choi (1995) show,

$$\forall m \in m_{\underline{F}_{\tilde{H}, \tilde{c}}}(\mathbb{C}^+) \quad z_{\underline{F}_{\tilde{H}, \tilde{c}}}(m) = -\frac{1}{m} + \frac{\tilde{c}}{m} - \frac{\tilde{c}}{m^2} m_{\tilde{H}} \left( -\frac{1}{m} \right) ,$$

which, letting  $v \equiv -1/m$ , is equivalent to

$$\forall v \in \mathbb{C}^+ \text{ such that } -\frac{1}{v} \in m_{\underline{F}_{\tilde{H}, \tilde{c}}}(\mathbb{C}^+) \quad m_{\tilde{H}}(v) = -\frac{1}{\tilde{c} v^2} \left[ z_{\underline{F}_{\tilde{H}, \tilde{c}}} \left( -\frac{1}{v} \right) - v + \tilde{c} v \right] . \quad (\text{A.16})$$

For the special case of  $\tilde{H} \equiv H$  and  $\tilde{c} \equiv c$ , this simplifies to

$$\forall v \in \mathbb{C}^+ \text{ such that } -\frac{1}{v} \in m_{\underline{F}}(\mathbb{C}^+) \quad m_H(v) = -\frac{1}{c v^2} \left[ z_{\underline{F}} \left( -\frac{1}{v} \right) - v + c v \right] . \quad (\text{A.17})$$

Let  $M \subseteq \mathbb{C}^+$  be a compact set contained in  $m_{\underline{F}}(\mathbb{C}^+)$  and also contained in  $m_{\underline{F}_{\hat{H}_n, \hat{c}_n}}(\mathbb{C}^+)$ , at least for  $n$  large enough. Let  $\{m_j\} \subseteq M$  be an infinite sequence with limit point  $m_0 \in M$ . Let  $v_j \equiv -1/m_j$  and  $v_0 \equiv -1/m_0$ . Then  $\{v_j\} \subseteq \mathbb{C}^+$  with limit point  $v_0 \in \mathbb{C}^+$ . Finally, let  $z_j \equiv z_{\underline{F}}(m_j)$  and  $z_0 \equiv z_{\underline{F}}(m_0)$ .

Part (i) of the theorem implies that  $\underline{F}_{\hat{H}_n, \hat{c}_n} \Rightarrow \underline{F}$  a.s. It then follows from Corollary 1 of Geronimo and Hill (2003) that

$$m_{\underline{F}_{\hat{H}_n, \hat{c}_n}}(z_j) \rightarrow m_{\underline{F}}(z_j) \text{ a.s. } \forall j .$$

In particular, the proof of Corollary 1 of Geronimo and Hill (2003) uses that convergence in distribution of probability measures implies convergence of integrals of bounded and continuous functions. A completely analogous argument can therefore be invoked to show that also

$$z_{\underline{F}_{\hat{H}_n, \hat{c}_n}}(m_j) \rightarrow z_{\underline{F}}(m_j) \quad \text{a.s. } \forall j$$

or, equivalently, that

$$z_{\underline{F}_{\hat{H}_n, \hat{c}_n}}\left(-\frac{1}{v_j}\right) \rightarrow z_{\underline{F}}\left(-\frac{1}{v_j}\right) \quad \text{a.s. } \forall j .$$

Using relation (A.16), with  $\tilde{H} \equiv \hat{H}_n$  and  $\tilde{c} \equiv \hat{c}_n$ , and relation (A.17), this implies that

$$\begin{aligned} m_{\hat{H}_n}(v_j) &= -\frac{1}{\hat{c}_n v_j^2} \left[ z_{\underline{F}_{\hat{H}_n, \hat{c}_n}}\left(-\frac{1}{v_j}\right) - v_j + \hat{c}_n v_j \right] \\ &\rightarrow -\frac{1}{c v_j^2} \left[ z_{\underline{F}}\left(-\frac{1}{v_j}\right) - v_j + c v_j \right] = m_H(v_j) \quad \text{a.s. } \forall j , \end{aligned}$$

which completes the proof of part (ii) the theorem. ■

**PROOF OF COROLLARY 5.1.** We start with the proof of part (i). Following El Karoui (2008), we call  $H_{T_n}$  a discretization of  $H$  on the grid  $\{J_n/T_n, (J_n + 1)/T_n, \dots, K_n/T_n\}$ . For instance, we can choose  $H_{T_n}$  to be a step function with  $H_{T_n}(x) \equiv H(x)$  if  $x = l/T_n$ ,  $l \in \mathbb{N}$ , and  $H_{T_n}$  is constant on  $[l/T_n, (l + 1)/T_n)$ . If the support of  $H$  is given by  $[h_1, h_2]$ , say, then the support of  $H_{T_n}$  is contained in  $[h_1 - 1/T_n, h_2 + 1/T_n]$ . It is easy to see that for such a discretization  $H_{T_n}$ , it holds that  $H_{T_n} \Rightarrow H$ , as long as

$$\exists b > 0 \text{ such that } \lambda_p \leq b \text{ for all } n \text{ sufficiently large} \quad \text{and} \quad (\text{A.18})$$

$$\exists \gamma > 0 \text{ such that } J_n/T_n \leq h_1 - \gamma \text{ and } K_n/T_n \geq h_2 + \gamma \text{ for all } n \text{ sufficiently large} . \quad (\text{A.19})$$

First, (A.18) holds a.s. as shown by Bai and Silverstein (1998) and Mestre (2008, Section II), given our set of assumptions, in particular Assumption (A4). Second, the support of  $F$  is denoted by  $[\tilde{z}_1, \tilde{z}_2]$ . On the one hand, it follows from Lemma 1.4 of Bai and Silverstein (1999) that  $\tilde{z}_1 < h_1$  and  $\tilde{z}_2 > h_2$ . Therefore, it holds that  $z_1 = h_1 - \delta_1$  and  $z_2 = h_2 + \delta_2$  for some  $\delta_1, \delta_2 > 0$ . On the other hand,  $F_n \Rightarrow F$  a.s., implying that  $\lambda_1 \leq \tilde{z}_1 + \delta_1/2$  and  $\lambda_p \geq \tilde{z}_2 - \delta_2/2$  for  $n$  sufficiently large a.s. So, letting  $\gamma \equiv \min\{\delta_1/2, \delta_2/2\}$ , condition (A.19) holds a.s. as well. Taken together, it follows that  $H_{T_n} \Rightarrow H$  a.s.

By construction,

$$\|F_{\hat{H}_n, \hat{c}_n} - F_n\|_{Q_n} \leq \|F_{H_{T_n}, \hat{c}_n} - F_n\|_{Q_n} \leq \|F_{H_{T_n}, \hat{c}_n} - F\|_{Q_n} + \|F - F_n\|_{Q_n} .$$

We know that  $\|F - F_n\|_{Q_n} \rightarrow 0$  a.s. So to establish part (i), it is sufficient to show that  $\|F_{H_{T_n}, \hat{c}_n} - F\|_{Q_n} \rightarrow 0$  a.s. Since  $H_{T_n} \Rightarrow H$  a.s. and  $\hat{c}_n \rightarrow c$ , it follows from Lemma A.2 and Polya's Theorem that  $\sup_t |F_{H_{T_n}, \hat{c}_n}(t) - F(t)| \rightarrow 0$  a.s., implying that  $\|F_{H_{T_n}, \hat{c}_n} - F\|_{Q_n} \rightarrow 0$  a.s.

But, having established part (i), part (ii) follows in exactly the same fashion as in the proof of Theorem 5.1. ■

PROOF OF COROLLARY 5.2. We start with some preliminary results, leading up to the proof of part (ii). Let  $G$  be a c.d.f. with continuous density  $g$  and compact support  $[a, b]$ . For a grid  $Q \equiv \{\dots, t_{-1}, t_0, t_1, \dots\}$  covering the support of  $G$ , the approximation to  $G$  via trapezoidal integration is defined as in (5.3). Since  $g$  is Lipschitz-continuous on  $[a, b]$ , there exists a (smallest) finite  $\varepsilon > 0$  such that  $|g(t_1) - g(t_2)| \leq \varepsilon$  as long as  $|t_1 - t_2| \leq \gamma$ . Denote by  $\hat{g}_Q$  the density corresponding to  $\hat{G}_Q$ . By definition of the trapezoidal rule,  $\hat{g}_Q$  is piecewise linear and agrees with  $g$  at all points  $t_k \in Q$ . Since the grid size of  $Q$  is given by  $\gamma$ , we may infer that

$$\sup_t |g(t) - \hat{g}_Q(t)| \leq 2\varepsilon \quad \text{and thus} \quad \sup_t |G(t) - \hat{G}_Q(t)| \leq 2\varepsilon(b - a + 2\gamma). \quad (\text{A.20})$$

We have assumed from the outset that  $c < 1$ . By construction,

$$\|\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F_n\|_{Q_n} \leq \|\hat{F}_{H_{T_n}, \hat{c}_n; Q_n} - F_n\|_{Q_n} \leq \|\hat{F}_{H_{T_n}, \hat{c}_n; Q_n} - F_{H_{T_n}, \hat{c}_n}\|_{Q_n} + \|F_{H_{T_n}, \hat{c}_n} - F_n\|_{Q_n}.$$

It follows from the proof of Corollary 5.1 that  $\|F_{H_{T_n}, \hat{c}_n} - F_n\|_{Q_n} \rightarrow 0$  a.s. So if we can show that  $\|\hat{F}_{H_{T_n}, \hat{c}_n; Q_n} - F_{H_{T_n}, \hat{c}_n}\|_{Q_n} \rightarrow 0$ , it follows that  $\|\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F_n\|_{Q_n} \rightarrow 0$  a.s.

For any probability measure  $\tilde{H}$ , any  $\tilde{c} > 0$ , and any  $\lambda \in (0, +\infty)$ , let

$$\check{m}_{F_{\tilde{H}}, \tilde{c}}(\lambda) = \lim_{z \in \mathbb{C}^+ \rightarrow \lambda} m_{F_{\tilde{H}}, \tilde{c}}(z).$$

Also let  $f_{\tilde{H}, \tilde{c}}(\lambda) \equiv \pi^{-1} \text{Im}[\check{m}_{F_{\tilde{H}}, \tilde{c}}(\lambda)]$  and define  $f_{\tilde{H}, \tilde{c}}(0) \equiv 0$ . Then

$$\int_{-\infty}^t f_{\tilde{H}, \tilde{c}}(\lambda) d\lambda = \begin{cases} F_{\tilde{H}, \tilde{c}}(t) & \text{if } \tilde{c} < 1 \\ \tilde{c} \underline{F}_{\tilde{H}, \tilde{c}}(t) & \text{if } \tilde{c} > 1 \end{cases}.$$

We know that  $f \equiv f_{H, c}$  is continuous, and therefore Lipschitz-continuous, on  $[\tilde{z}_1, \tilde{z}_2]$  and constantly equal to zero outside  $[\tilde{z}_1, \tilde{z}_2]$ . Denote by  $f_{\max}$  the maximum value of  $f$ . Since  $H_{T_n} \Rightarrow H$ , it follows from part (i) of Proposition 4.2 that, for every  $\tilde{\delta} > 0$ ,

$$f_{H_{T_n}, \hat{c}_n}(\lambda) \rightarrow f(\lambda) \quad \text{uniformly in } \lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]. \quad (\text{A.21})$$

In particular, for every  $\varepsilon > 0$ , we can find  $N$  such that, for all  $n \geq N$ ,

$$|f_{H_{T_n}, \hat{c}_n}(\lambda) - f(\lambda)| < \varepsilon \quad \text{for all } \lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}].$$

For every  $n$ , the function  $f_{H_{T_n}, \hat{c}_n}$  is monotonically increasing near the left boundary of its support and monotonically decreasing near the right boundary of its support; see Silverstein and Choi (1995, Section 5). The compact support of  $F$  is given by  $[\tilde{z}_1, \tilde{z}_2]$ . Lemma A.2 then implies that the support of  $F_{H_{T_n}, \hat{c}_n}$  is contained in  $[\tilde{z}_1 - \eta_n, \tilde{z}_2 + \eta_n]$  for some positive sequence  $\eta_n \rightarrow 0$ , so

$$f_{H_{T_n}, \hat{c}_n}(\lambda) = 0 \quad \text{for } \lambda \notin [\tilde{z}_1 - \eta_n, \tilde{z}_2 + \eta_n]. \quad (\text{A.22})$$

And further, for  $\eta_n$  and  $\tilde{\delta}$  sufficiently small and for  $n$  sufficiently large, we may assume that

$$f_{H_{T_n}, \hat{c}_n}(\lambda) \leq 2f_{\max} \quad \text{for all } \lambda \in [\tilde{z}_1 - \eta_n, \tilde{z}_1 + \tilde{\gamma}_n] \cup [\tilde{z}_2 - \tilde{\gamma}_n, \tilde{z}_2 + \eta_n]. \quad (\text{A.23})$$

Since  $f$  is Lipschitz-continuous on  $[\tilde{z}_1, \tilde{z}_2]$ , for  $\varepsilon > 0$ , there exists  $\delta > 0$  such that  $|f(\lambda_1) - f(\lambda_2)| \leq \varepsilon/2$  for all  $\lambda_1, \lambda_2 \in [\tilde{z}_1, \tilde{z}_2]$  with  $|\lambda_1 - \lambda_2| < \delta$ . From (A.21) it then follows that for  $n$  large enough,

$$|f_{H_{T_n}, \hat{c}_n}(\lambda_1) - f_{H_{T_n}, \hat{c}_n}(\lambda_2)| \leq \varepsilon \quad \text{for all } \lambda_1, \lambda_2 \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}] \text{ with } |\lambda_1 - \lambda_2| \leq \delta.$$

Applying the previous discussion for a general c.d.f.  $G$  and a general grid  $Q$  leading to (A.20) to the special cases of  $F_{H_{T_n}, \hat{c}_n}$  and  $Q_n$ , respectively, we thus obtain that, for  $n$  large enough (in particular, satisfying  $\gamma_n \leq \delta$ ),

$$\sup_{\lambda \in [\tilde{z}_1 + \tilde{\delta}, \tilde{z}_2 - \tilde{\delta}]} |f_{H_{T_n}, \hat{c}_n}(\lambda) - \hat{f}_{H_{T_n}, \hat{c}_n; Q_n}(\lambda)| \leq 2\varepsilon. \quad (\text{A.24})$$

Combining (A.22)–(A.24) yields, for  $\varepsilon$  and  $\tilde{\delta}$  small enough and for  $n$  large enough,

$$\sup_{\lambda \in \mathbb{R}} |F_{H_{T_n}, \hat{c}_n}(\lambda) - \hat{F}_{H_{T_n}, \hat{c}_n; Q_n}(\lambda)| \leq 2\varepsilon(\tilde{z}_2 - \tilde{z}_1 + 2\delta) + 4f_{\max}(\eta_n + \tilde{\delta}). \quad (\text{A.25})$$

Since the right hand side of (A.25) can be made arbitrarily small, we have established that  $\|\hat{F}_{H_{T_n}, \hat{c}_n; Q_n} - F_{H_{T_n}, \hat{c}_n}\|_{Q_n} \rightarrow 0$ , which implies that  $\|\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F_n\|_{Q_n} \rightarrow 0$  a.s., which in turn implies that

$$\|\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F\|_{Q_n} \leq \|\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F_n\|_{Q_n} + \|F_n - F\|_{Q_n} \rightarrow 0 \text{ a.s.} \quad (\text{A.26})$$

Lemma A.1 then tells us that  $\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} \Rightarrow F$  a.s.

We now show that this implies part (ii) of the corollary, namely that  $\hat{H}_n \Rightarrow H$  a.s. by means of contradiction. To this end, assume that  $\hat{H}_n \Rightarrow H$  a.s. is not the case. The sequence  $\{\hat{H}_n\}$  is tight a.s. This is because the upper bound of the support of  $H_n$  is given by  $K_n/T_n$  which, by definition of  $K_n$  satisfies  $K_n/T_n \leq \lambda_p + 1/T_n$ ; and we know from Bai and Silverstein (1998) that for any  $\varepsilon > 0$ ,  $\lambda_p \leq \tilde{z}_2 + \varepsilon$  for  $n$  large enough a.s. Similar for the lower bound, or simply use zero as very crude lower bound. Therefore, if  $\hat{H}_n \Rightarrow H$  a.s. is not the case, there then exists a probability measure  $H' \neq H$  and a subsequence  $\{n_k\}$  such that on a set with positive probability, we have  $\hat{H}_{n_k} \Rightarrow H'$ .

Similarly to an argument used in the proof of part (i) of Corollary 5.1 — with  $\hat{H}_{n_k}$  and  $H'$  now playing the roles of  $H_{T_n}$  and  $H$ , respectively — it then follows that  $\|F_{\hat{H}_{n_k}, \hat{c}_{n_k}} - F_{H', c}\|_{Q_{n_k}} \rightarrow 0$  on a set with positive probability. But it also holds that  $\|\hat{F}_{\hat{H}_{n_k}, \hat{c}_{n_k}; Q_{n_k}} - F_{\hat{H}_{n_k}, \hat{c}_{n_k}}\|_{Q_{n_k}} \rightarrow 0$  similarly to an argument used above — with  $F_{\hat{H}_{n_k}, \hat{c}_{n_k}}$  now playing the role of  $F_{H_{T_n}, \hat{c}_n}$ . Together, we obtain that  $\|\hat{F}_{\hat{H}_{n_k}, \hat{c}_{n_k}; Q_{n_k}} - F_{H', c}\|_{Q_{n_k}} \rightarrow 0$  on a set with positive probability. Since we are working under the assumption that  $c < 1$ , both  $F_H$  and  $F_{H'}$  are continuous. Lemma A.1 then tells us that  $\sup_t |\hat{F}_{\hat{H}_{n_k}, \hat{c}_{n_k}; Q_{n_k}}(t) - F_{H', c}(t)| \rightarrow 0$  on a set with positive probability. But this in contradiction to  $\sup_t |\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F(t)| \rightarrow 0$  a.s. So the proof of part (ii) is accomplished.

We now can establish that  $\|\hat{F}_{\hat{H}_n, \hat{c}_n; Q_n} - F_{\hat{H}_n, \hat{c}_n}\|_{Q_n} \rightarrow 0$  a.s., knowing that  $\hat{H}_n \Rightarrow H$  a.s., very much in the same way as we established before that  $\|\hat{F}_{H_{T_n}, \hat{c}_n; Q_n} - F_{H_{T_n}, \hat{c}_n}\|_{Q_n}$ , knowing

that  $H_{T_n} \Rightarrow H$ . As a result, we obtain that  $\|F_{\hat{H}_n, \hat{c}_n} - F\|_{Q_n} \rightarrow 0$  a.s. Invoking Lemma A.1 establishes part (i) then.

Parts (iii)–(iv) follow immediately from parts (i)–(ii) and Proposition 4.3, part (ii). ■

## B Additional Monte Carlo Simulations

### B.1 Comparisons with Other Estimators

So far, we have compared the nonlinear shrinkage estimator  $\hat{S}_n$  only to the linear shrinkage estimator  $\bar{S}_n$  and the oracle estimator  $S_n^{or}$  to keep the resulting figures concise and legible.

It is of additional interest to compare the nonlinear shrinkage estimator also to some other estimators from the literature. To this end we consider the following set of estimators:

- The estimator of Stein (1975).
- The estimator of Haff (1980).
- The estimator recently proposed by Won et al. (2009). This estimator is based on a maximum likelihood approach, assuming normality, with an explicit constraint on the condition number of the covariance matrix. The resulting estimator turns out to be a nonlinear shrinkage estimator as well: all ‘small’ sample eigenvalues are brought up to a lower bound, all ‘large’ sample eigenvalues are brought down to an upper bound, and all ‘intermediate’ sample eigenvalues are left unchanged.

Therefore, the corresponding transformation from sample eigenvalues to shrunk eigenvalues is step-wise linear: first flat, then a 45-degree line, and then flat again. The upper and lower bounds are determined by the desired constraint on the condition number  $\kappa$ . If such an explicit constraint is not available from *a priori* information, a suitable constraint number  $\hat{\kappa}$  can be computed in a data-dependent fashion by a  $K$ -fold cross-validation method, which is the method we use.<sup>2</sup>

In particular, the cross-validation method selects  $\hat{\kappa}$  by optimizing over a finite grid  $\{\kappa_1, \kappa_2, \dots, \kappa_L\}$  that has to be supplied by the user. To this end we choose  $L = 10$  and the  $\kappa_l$  log-linearly spaced between 1 and  $\kappa(S_n)$ , for  $l = 1, \dots, L$ ; here  $\kappa(S_n)$  denotes the condition number of the sample covariance matrix. More precisely, for  $l = 1, \dots, L$ ,  $\kappa_l \equiv \exp(\omega_l)$ , where  $\{\omega_1, \omega_2, \dots, \omega_L\}$  is the equally-spaced grid with  $\omega_1 \equiv 0$  and  $\omega_L \equiv \log(\kappa(S_n))$ .

- The cross-validation version of the nonlinear shrinkage estimator  $\hat{S}_n$ ; see Remark 5.2.

We repeat the simulation exercises of Subsections 6.1–6.3, replacing the oracle estimator and the linear shrinkage estimator with the above set of other estimators. The respective PRIALs of the various estimators are plotted in Figures 11–13.

---

<sup>2</sup>We are grateful to Joong-Ho Won for supplying us with corresponding Matlab code.

One can see that the nonlinear shrinkage estimator  $\hat{S}_n$  outperforms all other estimators, with the cross-validation version of  $\hat{S}_n$  in second place, followed by the estimators of Stein (1975), Won et al. (2009), and Haff (1980).

## B.2 Comparisons Based on a Different Loss Function

So far, the PRIAL has been based on the loss function

$$L^{Fr}(\hat{\Sigma}_n, \Sigma_n) \equiv \|\hat{\Sigma}_n - \Sigma_n\|^2 .$$

It is of additional interest to add some comparisons based on a different loss function. To this end we use the scale-invariant loss function proposed by James and Stein (1961), namely

$$L^{JS}(\hat{\Sigma}_n, \Sigma_n) = \text{trace}(\hat{\Sigma}_n \Sigma_n^{-1}) - \log \det(\hat{\Sigma}_n \Sigma_n^{-1}) - p . \quad (\text{B.1})$$

We repeat the simulation exercises of Subsections 6.1–6.3, replacing  $L^{Fr}$  with  $L^{JS}$ . The respective PRIALs of  $S_n^{or}$ ,  $\hat{S}_n$ , and  $\bar{S}_n$  are plotted in Figures 14–16.

One can see that the results do not change much qualitatively. If anything, the comparisons are now even more favorable to the nonlinear shrinkage estimator, in particular when comparing Figure 5 to Figure 16.

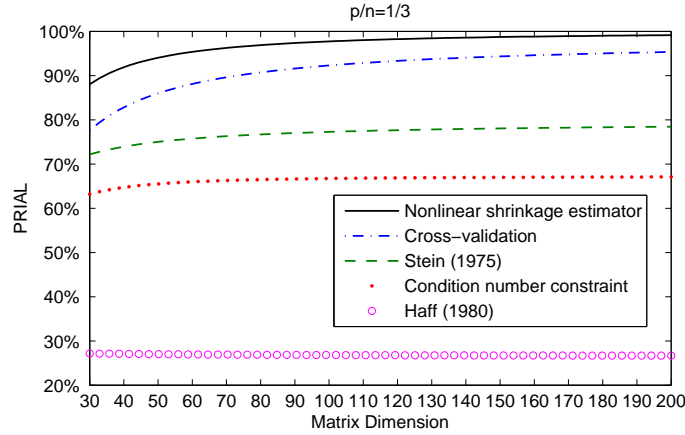


Figure 11: Comparison of Various Estimators. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. Every point is the result of 1,000 Monte Carlo simulations.

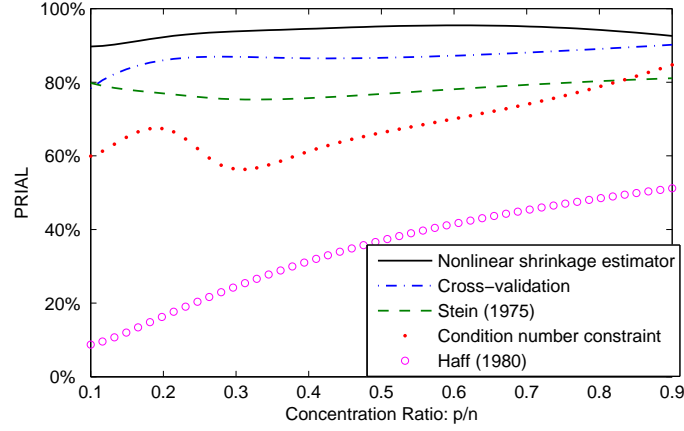


Figure 12: Effect of Varying the Concentration Ratio  $\hat{c}_n = p/n$ . 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. Every point is the result of 1,000 Monte Carlo simulations.

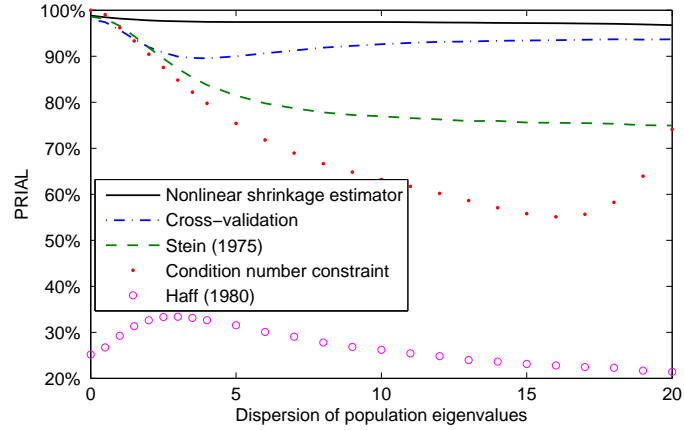


Figure 13: Effect of Varying the Dispersion of Population Eigenvalues. 20% of population eigenvalues are equal to 1, 40% equal to  $1 + 2d/9$ , and 40% equal to  $1 + d$ , where the dispersion parameter  $d$  varies from 0 to 20.  $p = 100$  and  $n = 300$ . Every point is the result of 1,000 Monte Carlo simulations.

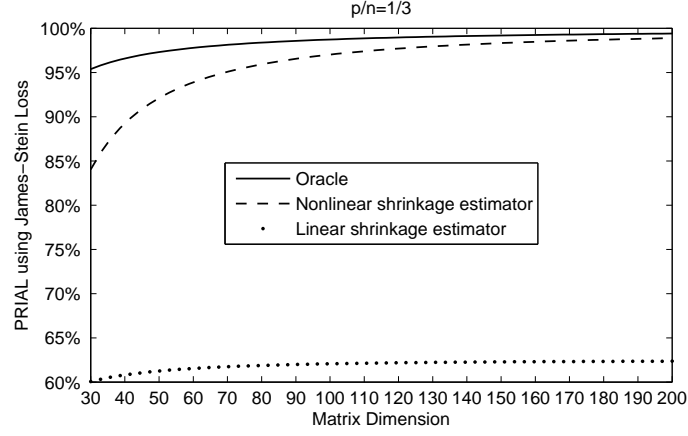


Figure 14: Comparison of the NonLinear vs. Linear Shrinkage Estimators. 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. The PRIALs are based on the James-Stein (1961) loss function (B.1). Every point is the result of 1,000 Monte Carlo simulations.

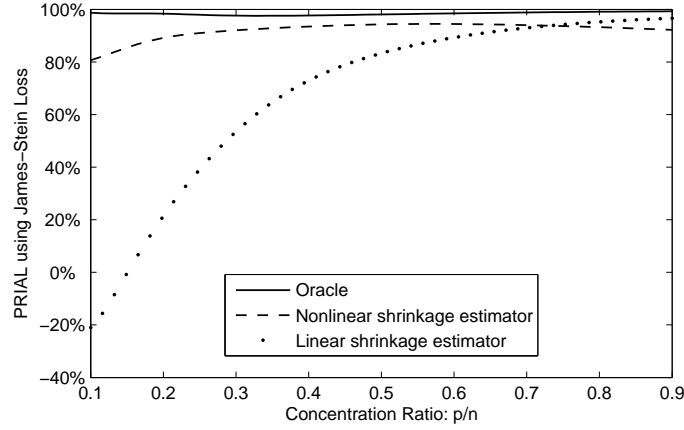


Figure 15: Effect of Varying the Concentration Ratio  $\hat{c}_n = p/n$ . 20% of population eigenvalues are equal to 1, 40% are equal to 3, and 40% are equal to 10. The PRIALs are based on the James-Stein (1961) loss function(B.1). Every point is the result of 1,000 Monte Carlo simulations.



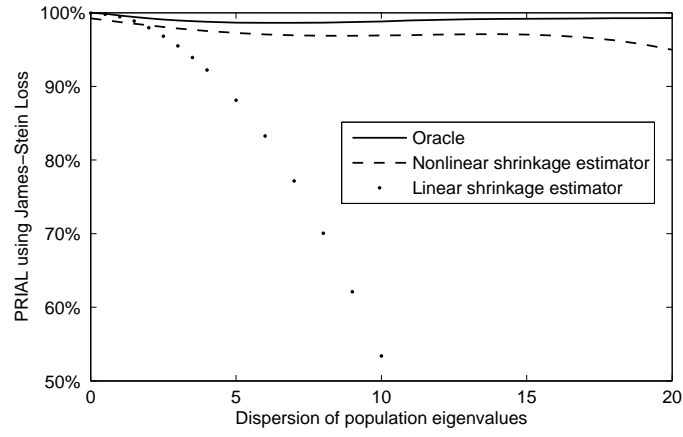


Figure 16: Effect of Varying the Dispersion of Population Eigenvalues. 20% of population eigenvalues are equal to 1, 40% equal to  $1 + 2d/9$ , and 40% equal to  $1 + d$ , where the dispersion parameter  $d$  varies from 0 to 20.  $p = 100$  and  $n = 300$ . The PRIALs are based on the James-Stein (1961) loss function (B.1). Every point is the result of 1,000 Monte Carlo simulations.